

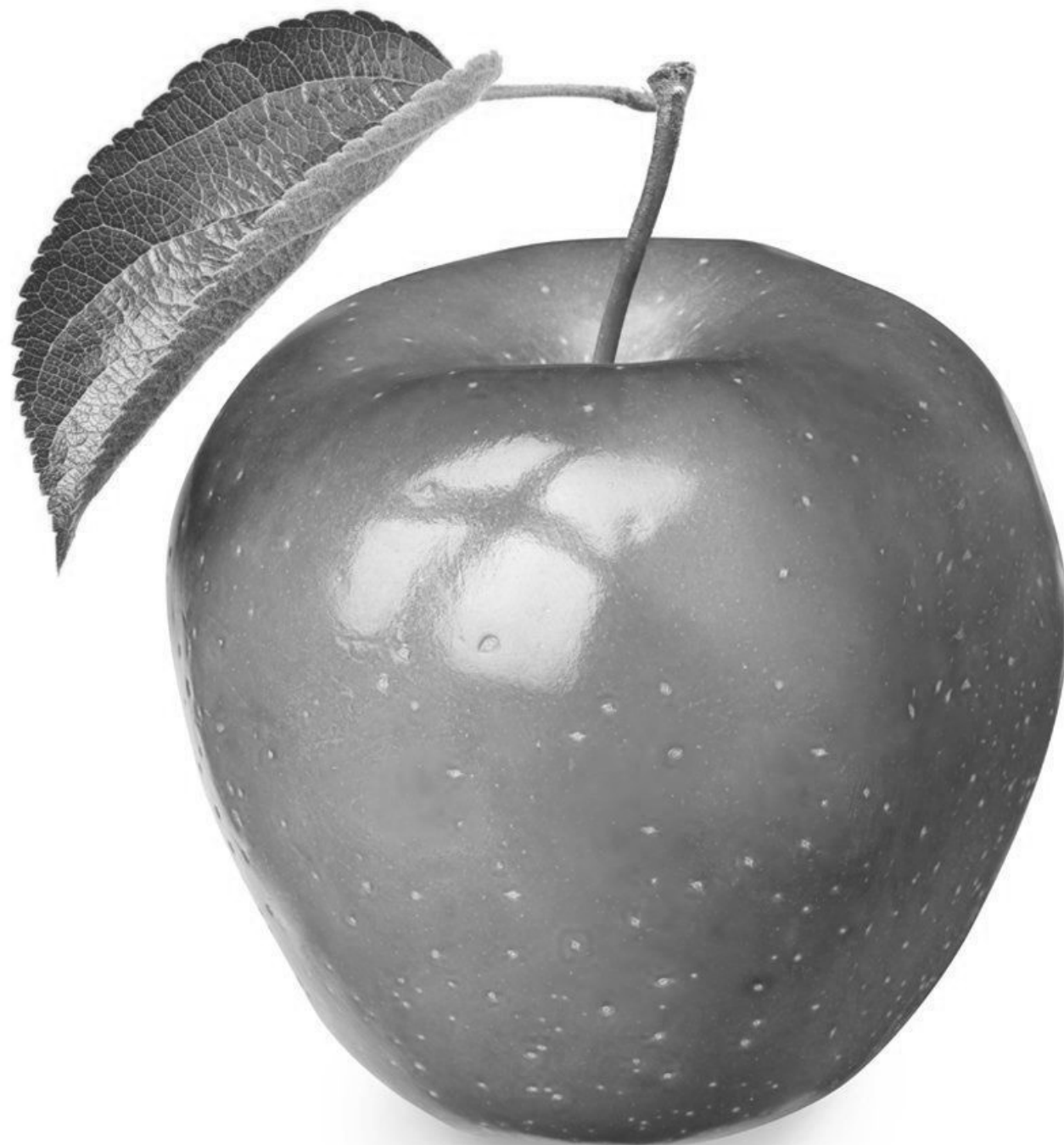
Vision and Language

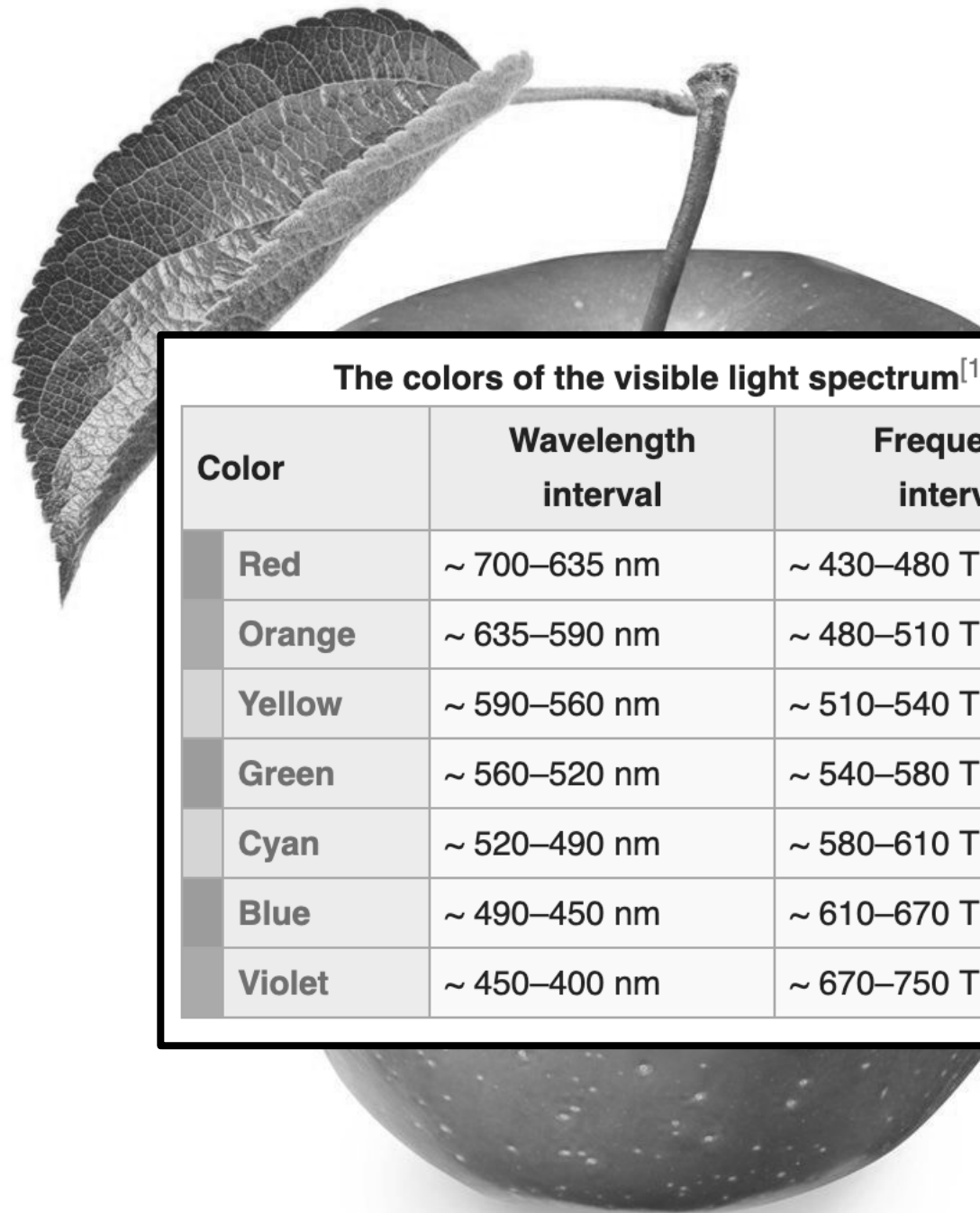


Rodolfo (Rudy) Corona

with thanks to Daniel Fried

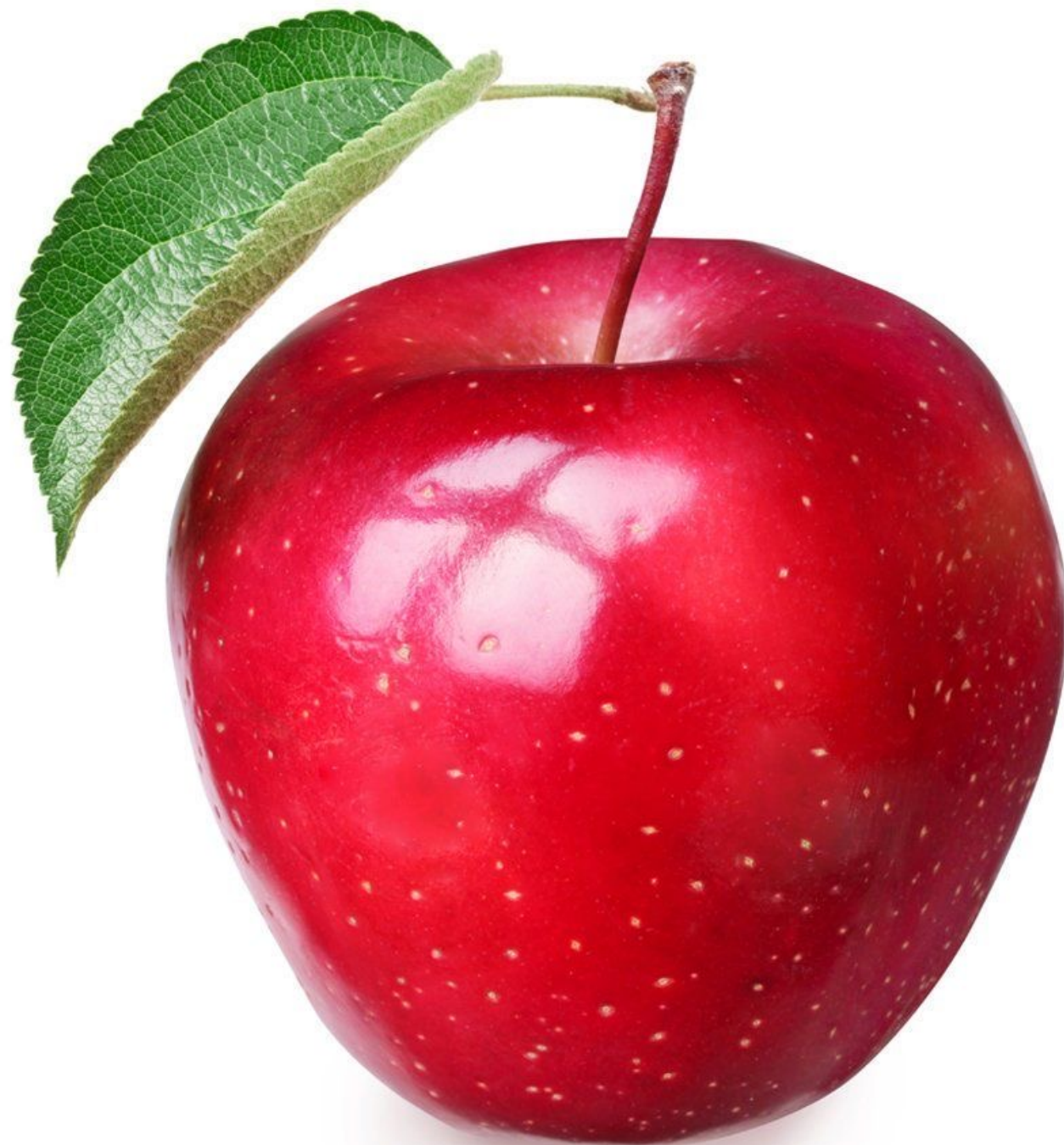
CS 288, 4/12/2022





The colors of the visible light spectrum^[1]

Color	Wavelength interval	Frequency interval
Red	~ 700–635 nm	~ 430–480 THz
Orange	~ 635–590 nm	~ 480–510 THz
Yellow	~ 590–560 nm	~ 510–540 THz
Green	~ 560–520 nm	~ 540–580 THz
Cyan	~ 520–490 nm	~ 580–610 THz
Blue	~ 490–450 nm	~ 610–670 THz
Violet	~ 450–400 nm	~ 670–750 THz



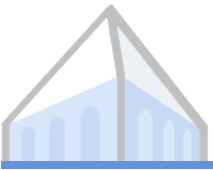
“Apples are red”

“The numbers this month are in the red”

“Red has a wavelength between 635-700nm”

...

“Pixel (1,1) has R=240, pixel (1,2) has ...”



What is Language Grounding?

- ▶ Tying language to non-linguistic things (e.g. a database in semantic parsing)
- ▶ The world only looks like a database some of the time!
- ▶ Some settings depend on grounding into *perceptual* or *physical* environments:

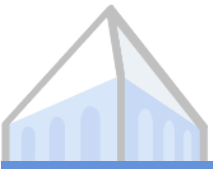


“Add the tomatoes and mix”



“Take me to the shop on the corner”

- ▶ **Focus today:** Grounding language to *visual perception*.



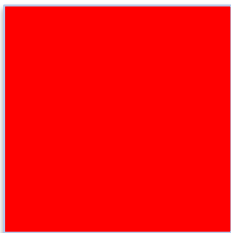
Grounding

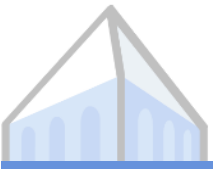
- (Some) possible things to ground into:



Grounding

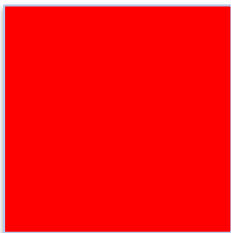
- ▶ (Some) possible things to ground into:
 - **Low-level percepts:** *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor...





Grounding

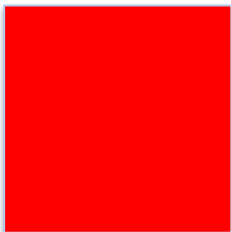
- ▶ (Some) possible things to ground into:
 - **Low-level percepts:** *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor...
 - **High-level percepts:** *cat* means this type of pattern





Grounding

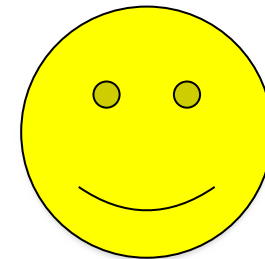
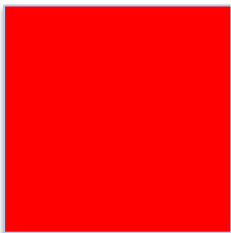
- ▶ (Some) possible things to ground into:
 - **Low-level percepts:** *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor...
 - **High-level percepts:** *cat* means this type of pattern
 - **Embodiment (effects on the world):** *go left* means the robot turns left, *speed up* means increasing actuation





Grounding

- ▶ (Some) possible things to ground into:
 - **Low-level percepts:** *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor...
 - **High-level percepts:** *cat* means this type of pattern
 - **Embodiment (effects on the world):** *go left* means the robot turns left, *speed up* means increasing actuation
 - **Social (effects on others):** polite language is correlated with longer forum discussions





Grounding

- ▶ (Some) possible things to ground into:
 - **Low-level percepts:** *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor...
 - **High-level percepts:** *cat* means this type of pattern
 - **Embodiment (effects on the world):** *go left* means the robot turns left, *speed up* means increasing actuation
 - **Social (effects on others):** polite language is correlated with longer forum discussions

For a nice taxonomy, related work, and examples, see *Experience Grounds Language* [Bisk et al. 2020]



Grounding

- ▶ (Some) key problems:
 - **Representation:** matching low-level percepts to high-level language (pixels vs *cat*)
 - **Abstraction and Composition:** meaning as a combination of parts
 - **Alignment:** aligning parts of language and parts of the world
 - **Content Selection and Context:** what are the important parts of the environment?
 - **Balance:** it's easy for multi-modal models to “cheat”, rely on imperfect heuristics, or ignore important parts of the input
 - **Generalization:** to novel world contexts / input combinations

The background image shows a large, light-colored stone building with a prominent clock tower. The clock tower has two visible clock faces and a belfry at the top. To the right, a portion of a classical building with large columns is visible. The sky is blue with some light clouds. A solid yellow horizontal bar is positioned below the text.

CS294-43: VISION AND LANGUAGE AI SEMINAR



A Gallery of Tasks

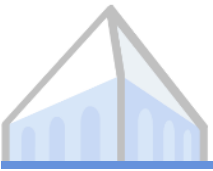


Image Captioning



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.



A horse carrying a large load of hay and two people sitting on it.



Bunk bed with a narrow shelf sitting underneath it.



Visual Question Answering

What is the dog wearing?
life jacket



collar



How many skiers are there?

2



1



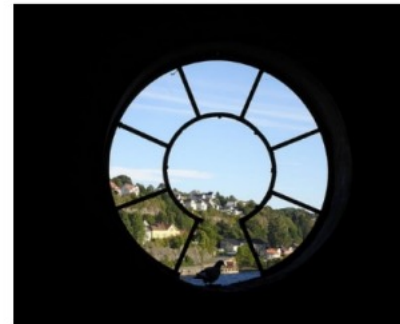
What number is on the train?
7907



8551



What is sitting in the window?
bird



clock





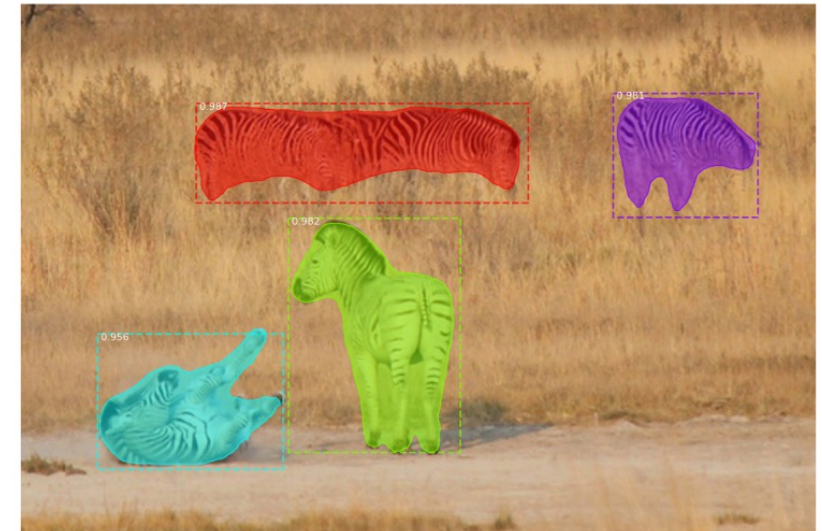
Object Detection (2D)



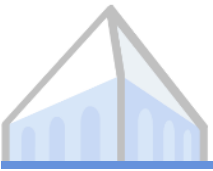
(a) Query: “street lamp”



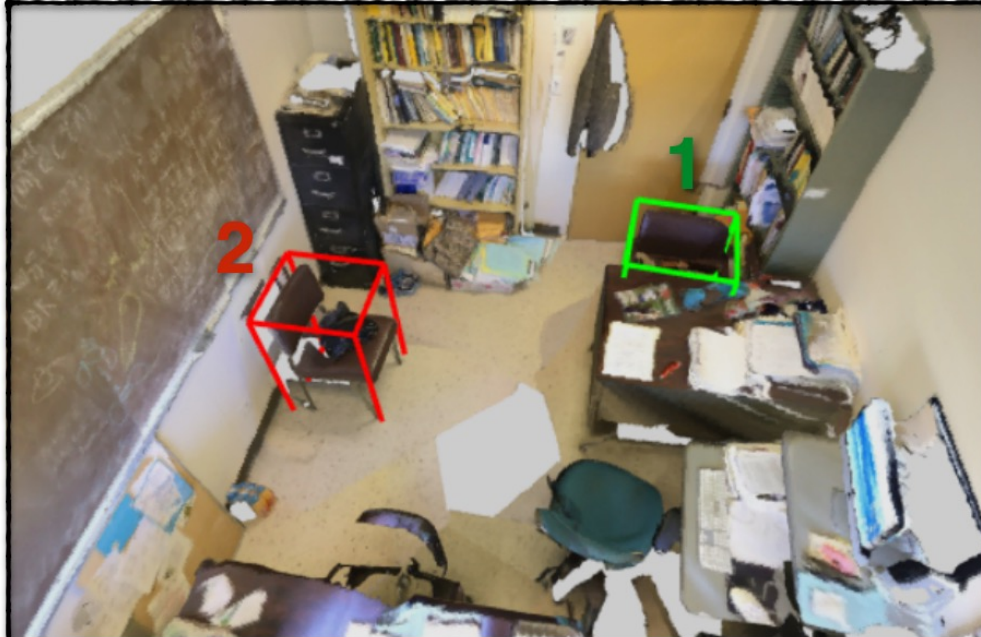
(b) Query: “major league logo”



(c) Query: “zebras on savanna”



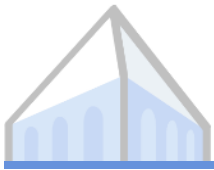
Object Detection (3D)



1. "The chair closest to the door."
2. "The chair under the chalkboard."



1. "The office chair that is green."
2. "Choose the brown office chair pushed under the desk."



Conditional Generation (2D)



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation



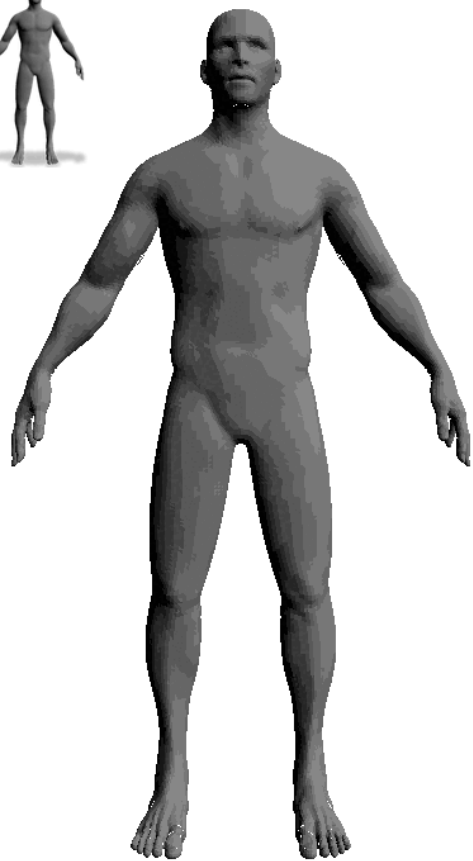
panda mad scientist mixing sparkling chemicals, artstation



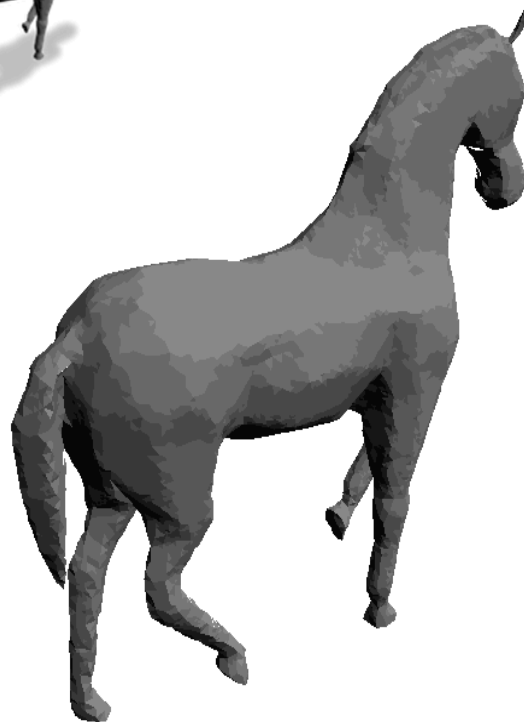
a corgi's head depicted as an explosion of a nebula



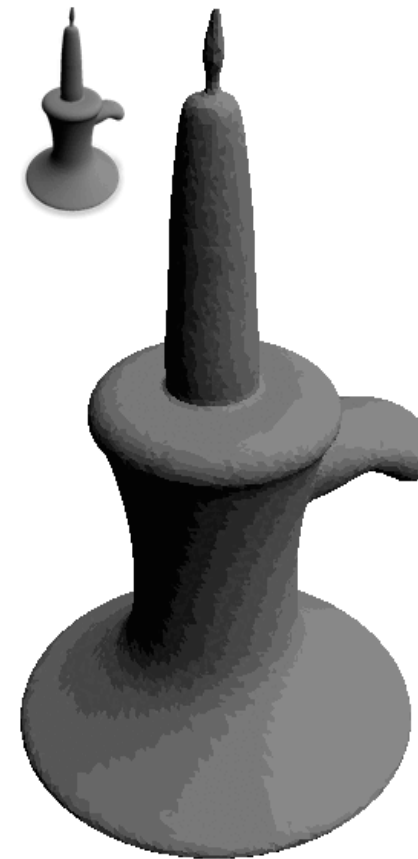
Conditional Generation (3D)



“Iron Man”



“Astronaut Horse”



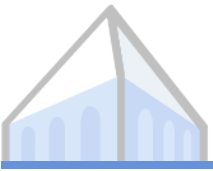
“Colorful Crochet Candle”



Vision and Language Navigation



“Place a clean ladle on a counter”



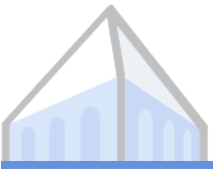
Why Grounded Language?

- Much language refers *to the world*.
- Convenient medium to communicate with machines!
- For many tasks, agents will need perceptual understanding and motor control for this interaction.

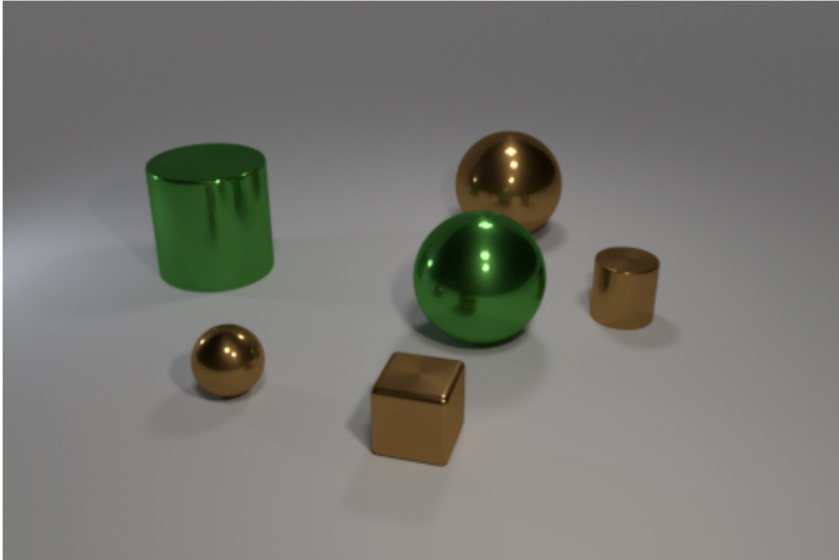




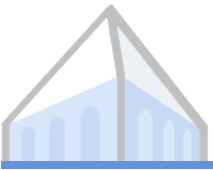




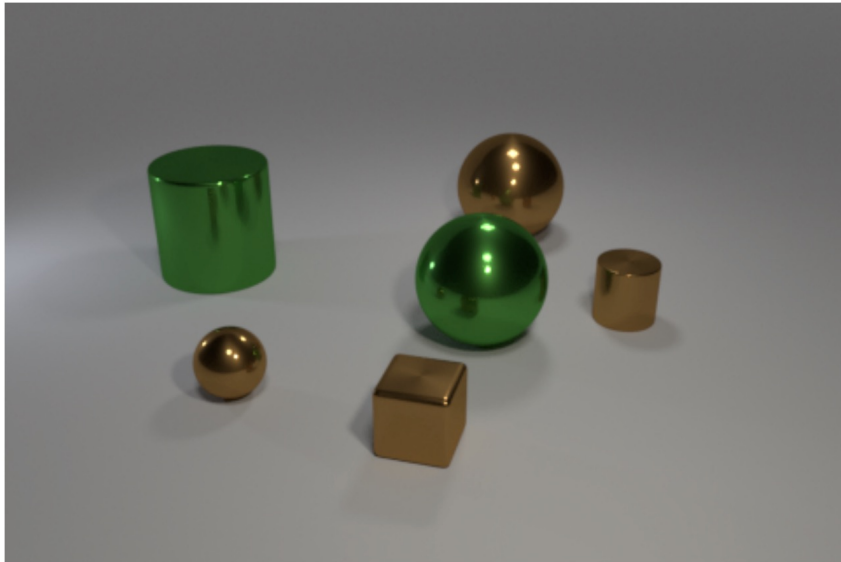
Bottom-Up & Top-Down Reasoning



“What color is the small
shiny cube?”



Bottom-Up & Top-Down Reasoning



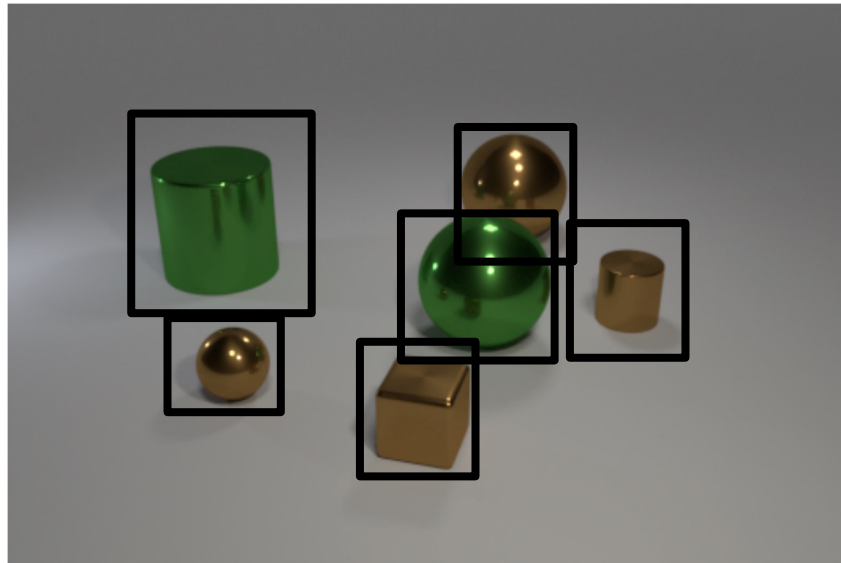
“What color is the small shiny cube?”

Neural
Network

ANSWER



Bottom-Up & Top-Down Reasoning

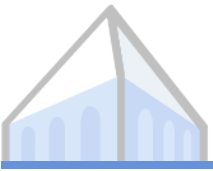


**Bottom-up
object proposals**

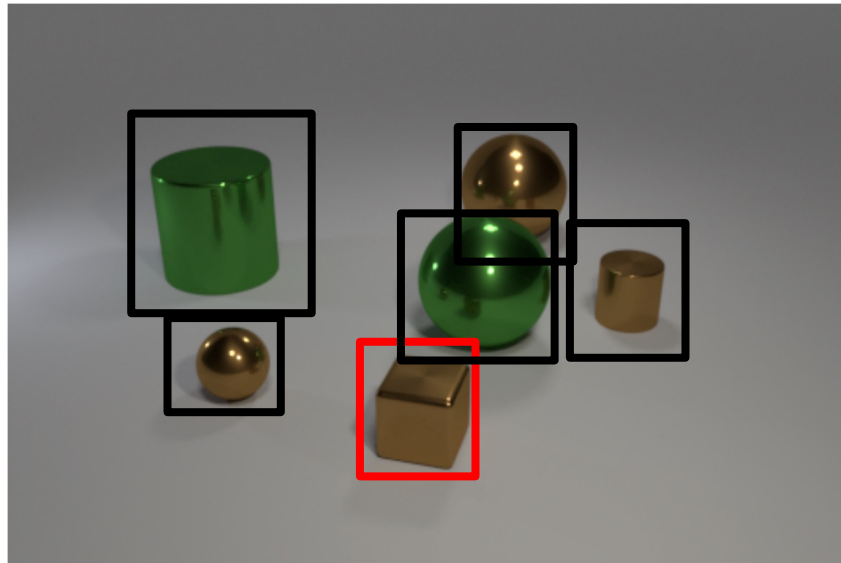
“What color is the small
shiny cube?”

al
Network

ANSWER



Bottom-Up & Top-Down Reasoning



**Top-down
attention**

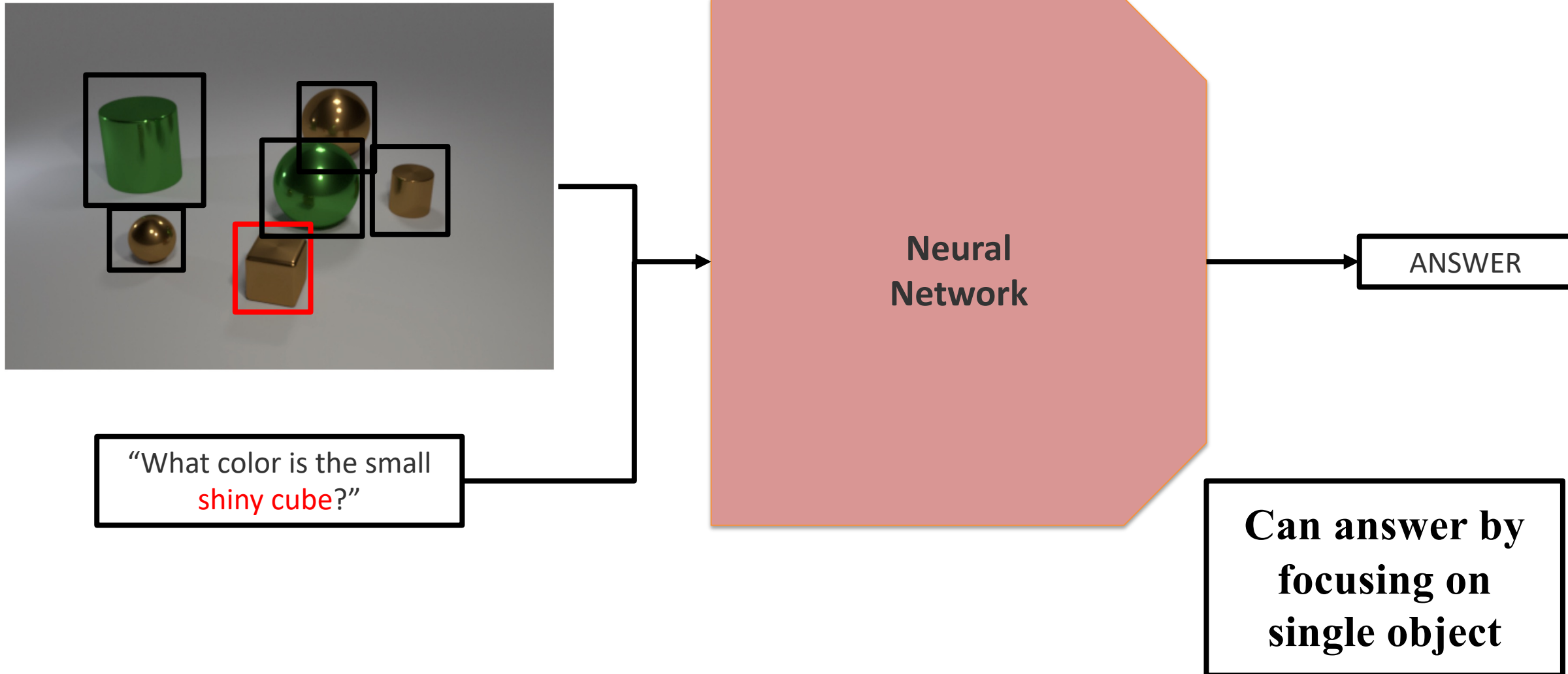
“What color is the small
shiny cube?”

al
Network

ANSWER



Bottom-Up & Top-Down Reasoning

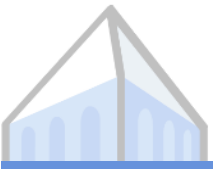




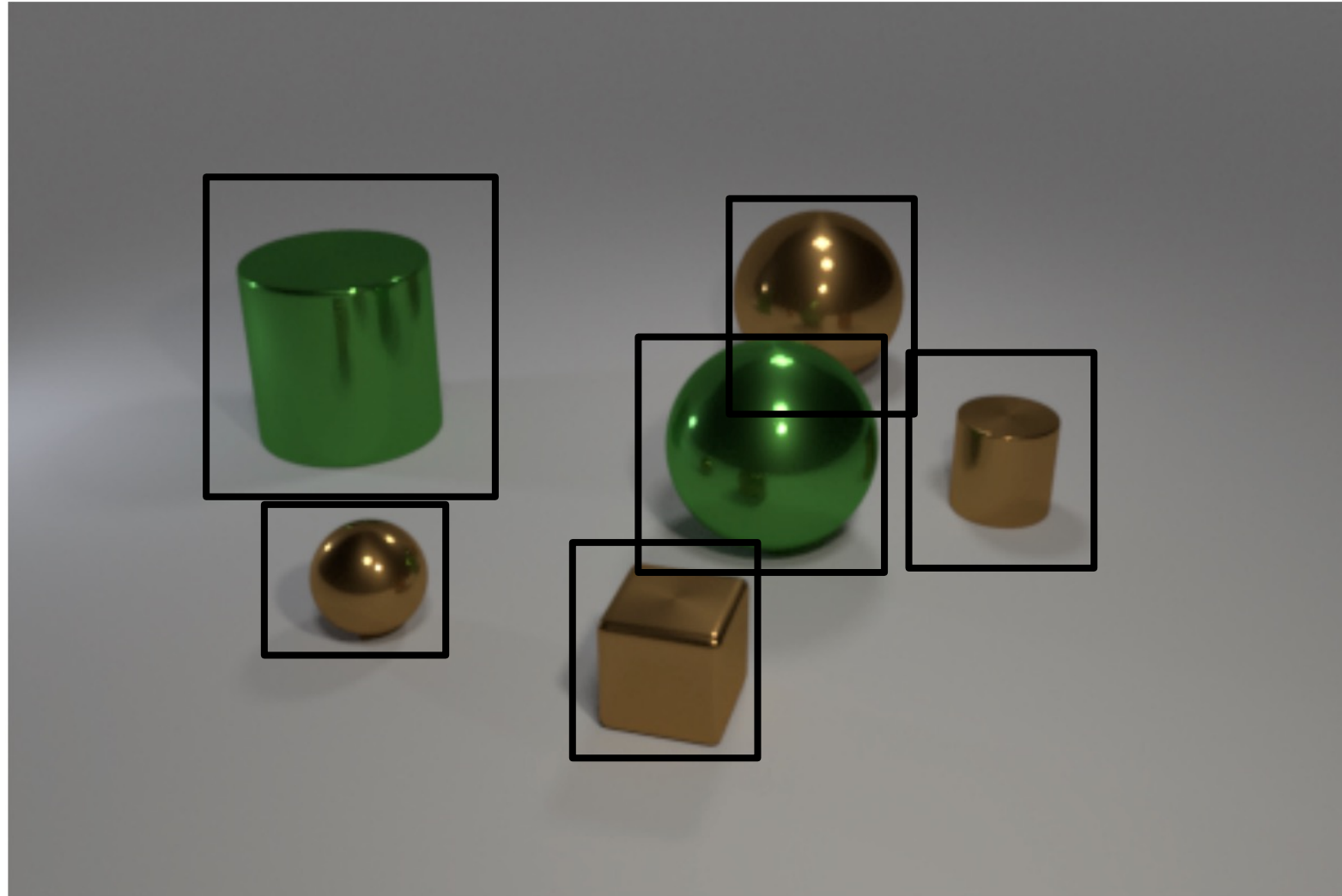
Bottom-Up & Top-Down Reasoning

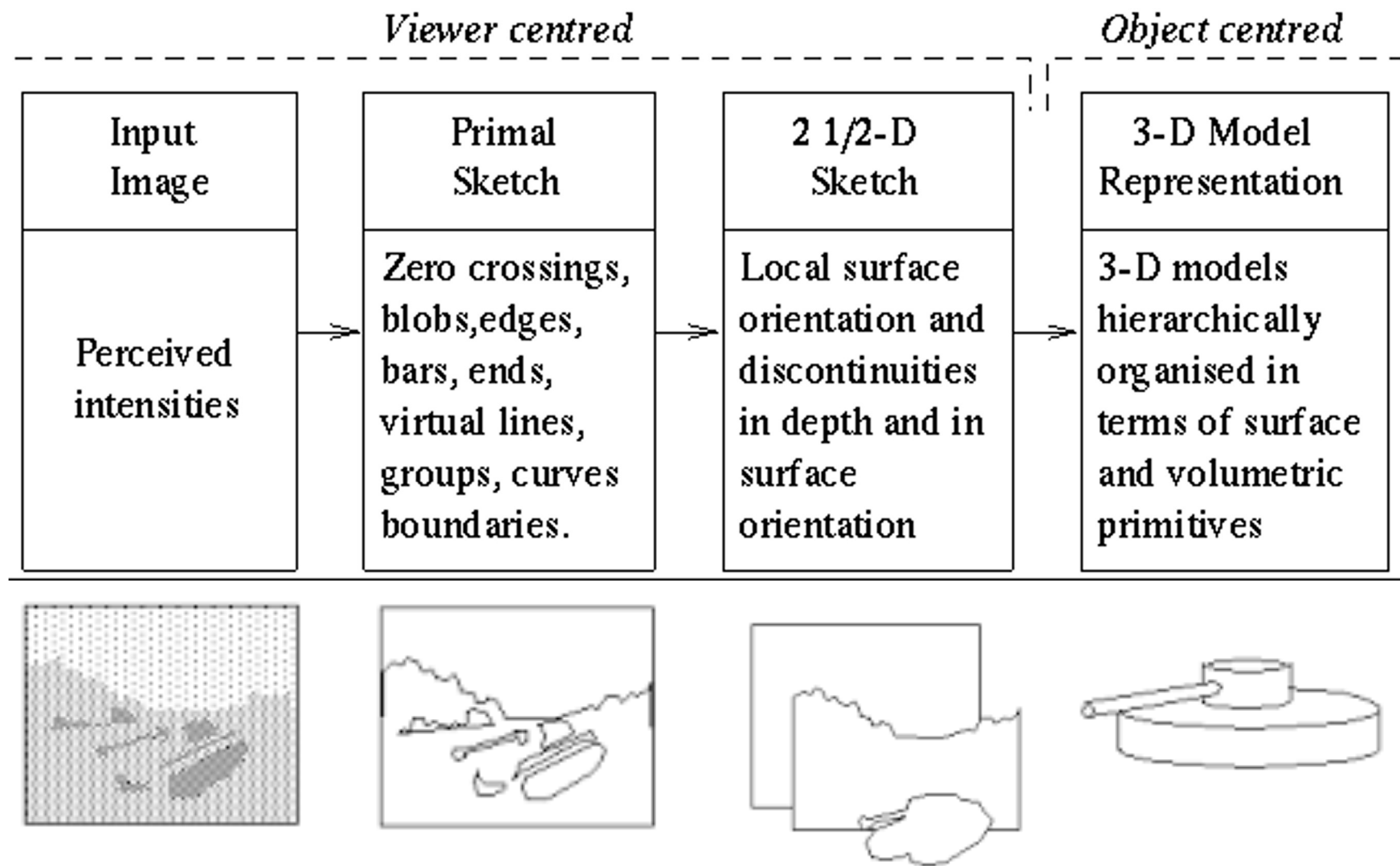
	Yes/No	Number	Other	Overall
Ours: ResNet (1×1)	76.0	36.5	46.8	56.3
Ours: ResNet (14×14)	76.6	36.2	49.5	57.9
Ours: ResNet (7×7)	77.6	37.7	51.5	59.4
Ours: Up-Down	80.3	42.8	55.8	63.2
Relative Improvement	3%	14%	8%	6%

**Provides inductive bias in
both directions!**



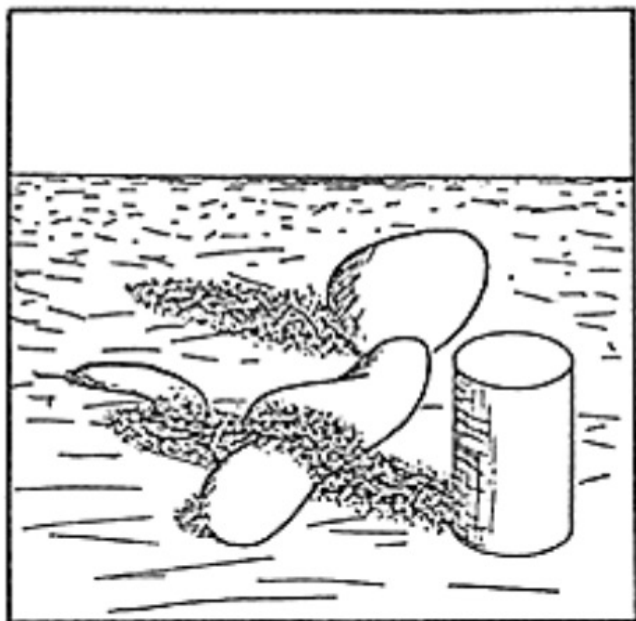
Bottom-Up



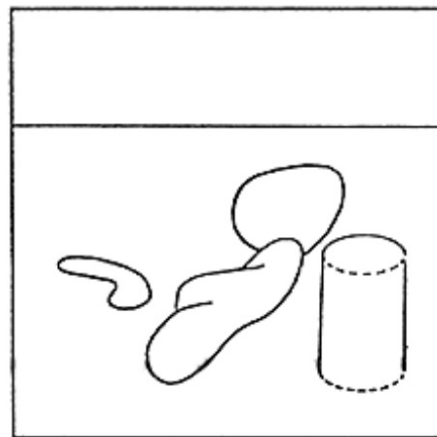




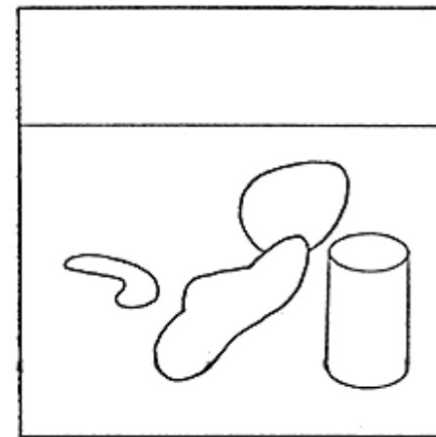
Intrinsic Images



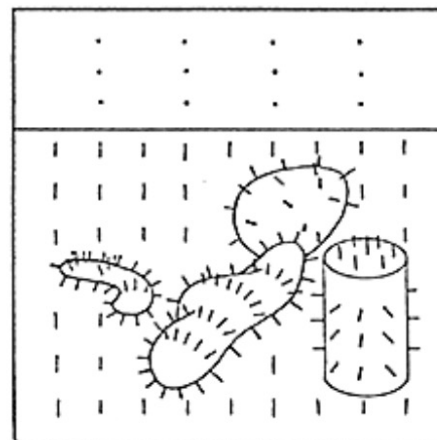
(a) ORIGINAL SCENE



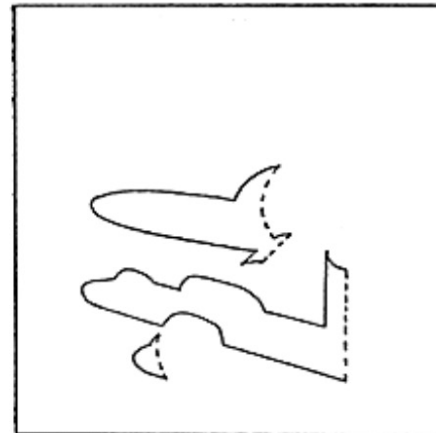
(b) DISTANCE



(c) REFLECTANCE



(d) ORIENTATION (VECTOR)



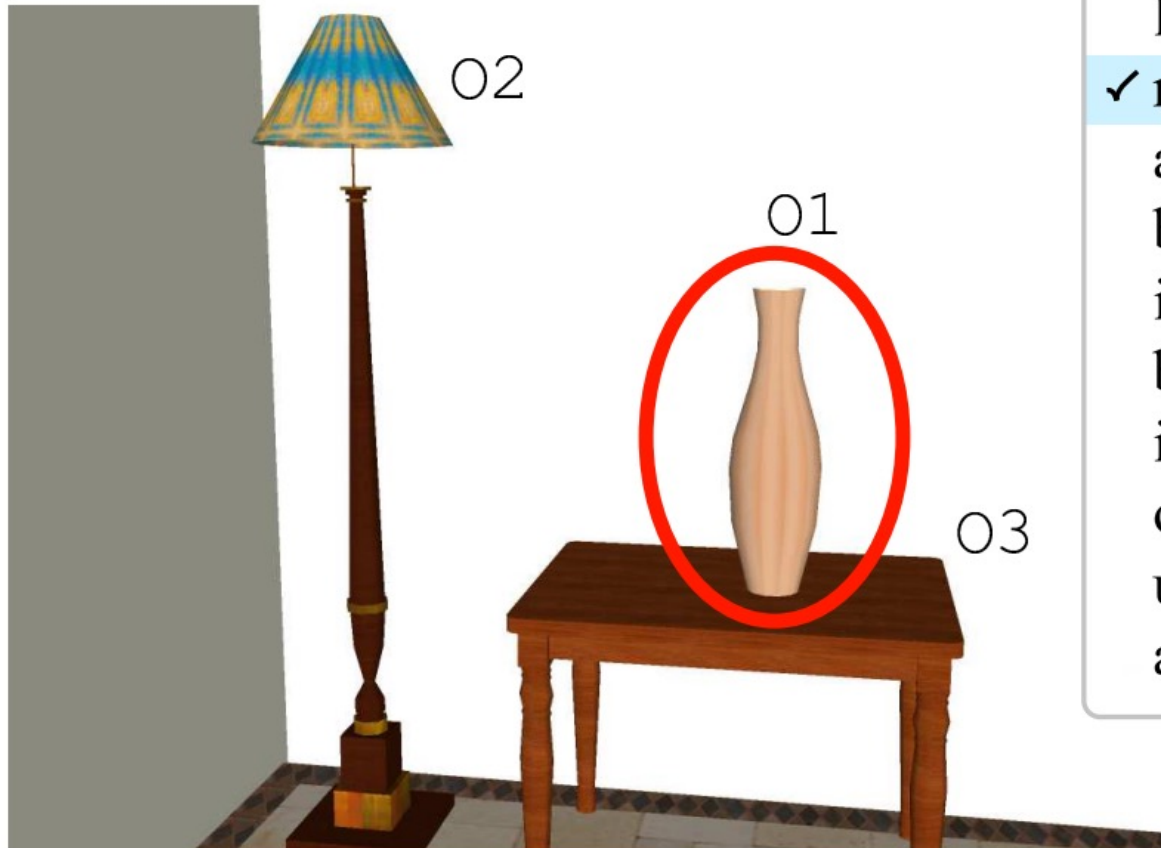
(e) ILLUMINATION



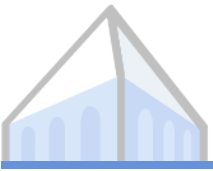
“Solved” Perception

Question: Where is the object outlined in red?

Answer: The object outlined in red is



left of
✓ right of
above
below
in front of
behind
inside of
on
under
across from

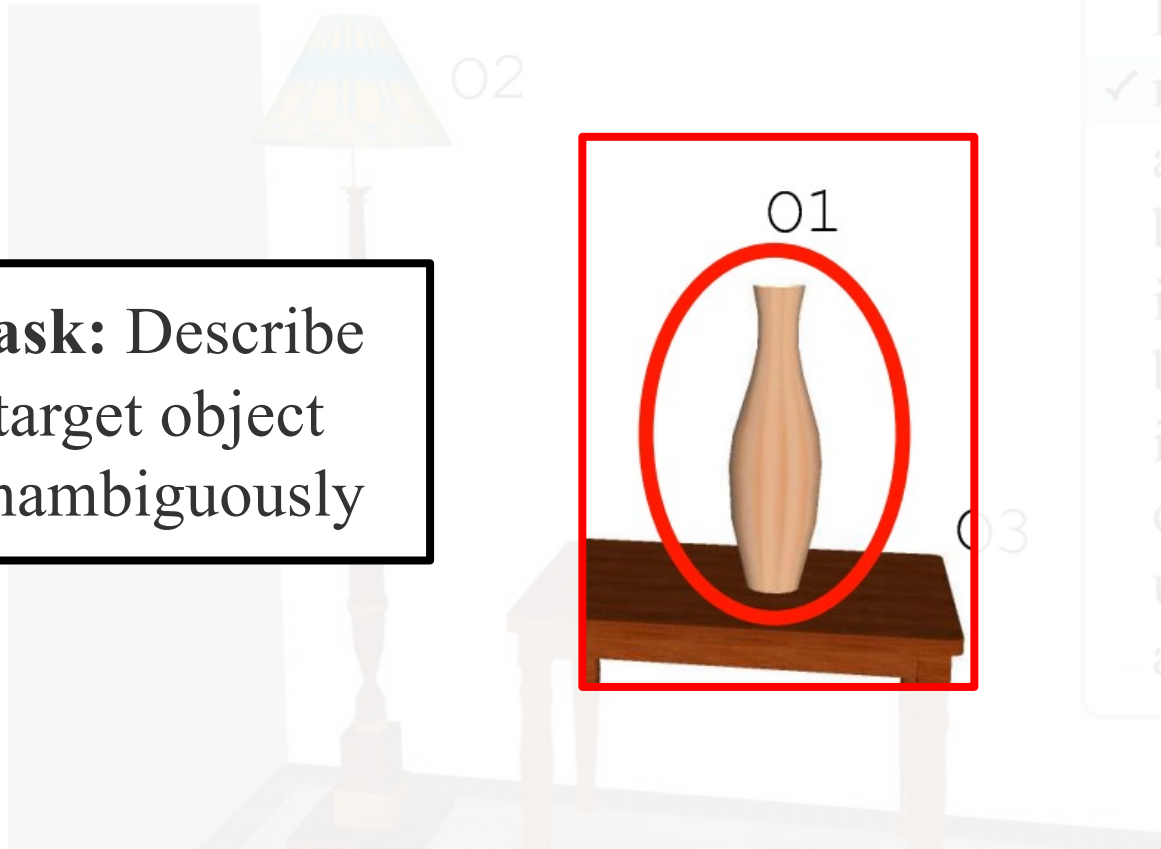


“Solved” Perception

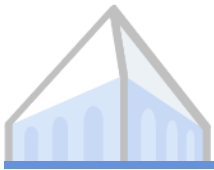
Question: Where is the object outlined in red?

Answer: The object outlined in red is

Task: Describe
target object
unambiguously



left of
✓ right of
above
below
in front of
behind
inside of
on
under
across from



“Solved” Perception

Question: Where is the object outlined in red?

Answer: The object outlined in red is

02

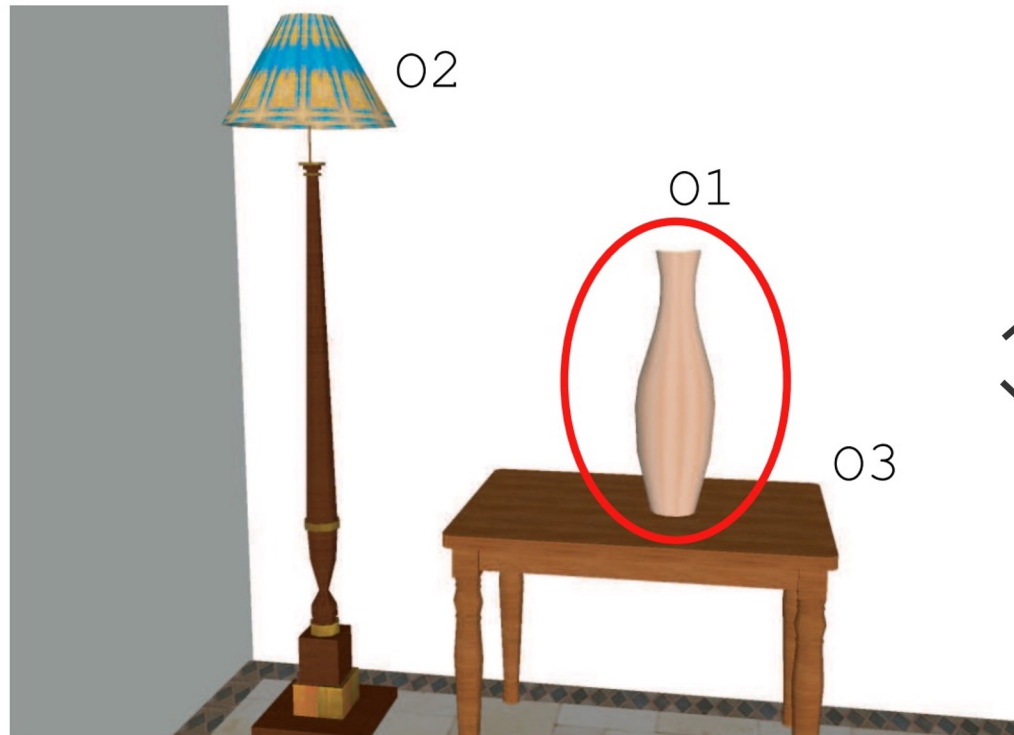
Relationships
between objects
known



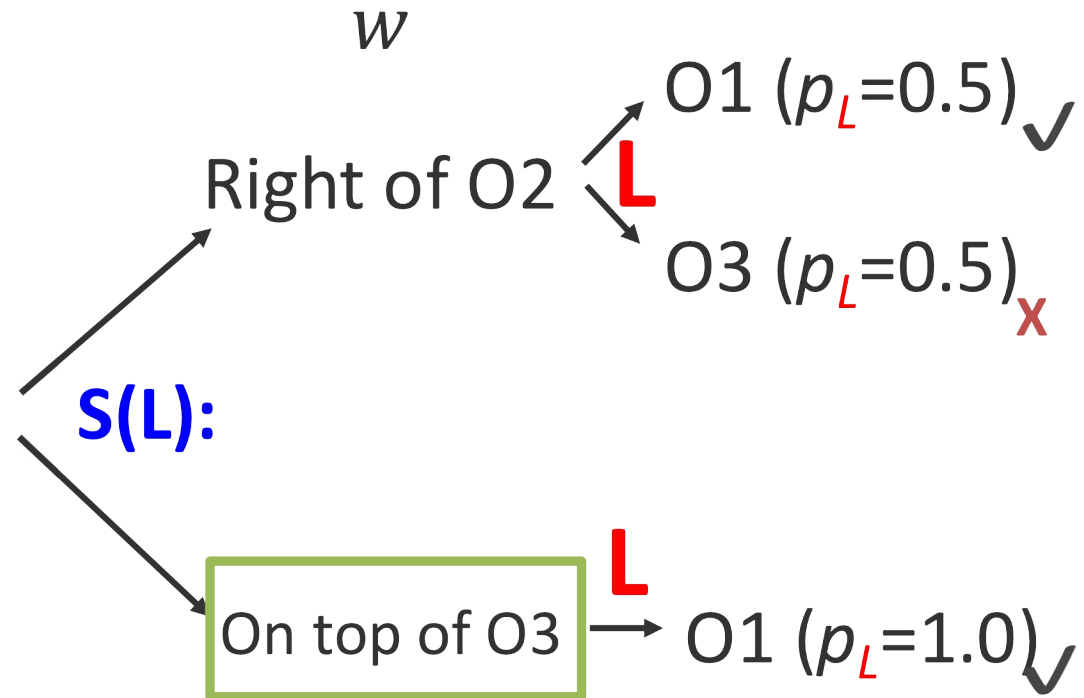
left of
✓ right of
above
below
in front of
behind
inside of
on
under
across from



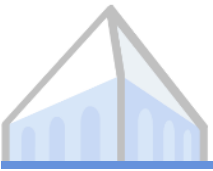
“Solved” Perception



Problem reduced to
pragmatic reasoning



$$S(L)(o) = \operatorname{argmax}_w p_L(o|w)$$



“Solved” Perception

“Go to the last butterfly on the right”



[(Cement, Easel, Cement, Butterfly, Wood, Butterfly),
(Wall, Empty, Wall, Butterfly, Wood, Butterfly),
(Cement, Empty, Wall, End, Wall, End)]



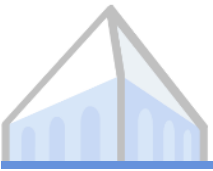
“Solved” Perception

“Go to the last butterfly on the right”

What annotators
see



[(Cement, Easel, Cement, Butterfly, Wood, Butterfly),
(Wall, Empty, Wall, Butterfly, Wood, Butterfly),
(Cement, Empty, Wall, End, Wall, End)]



“Solved” Perception

“Go to the last butterfly on the right”



What agent sees

[(Cement, Easel, Cement, Butterfly, Wood, Butterfly),
(Wall, Empty, Wall, Butterfly, Wood, Butterfly),
(Cement, Empty, Wall, End, Wall, End)]



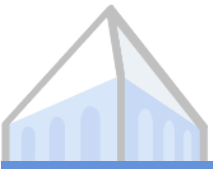
“Solved” Perception

“Go to the last butterfly on the right”

Reduced to
structured
prediction problem



[(Cement, Easel, Cement, Butterfly, Wood, Butterfly),
(Wall, Empty, Wall, Butterfly, Wood, Butterfly),
(Cement, Empty, Wall, End, Wall, End)]

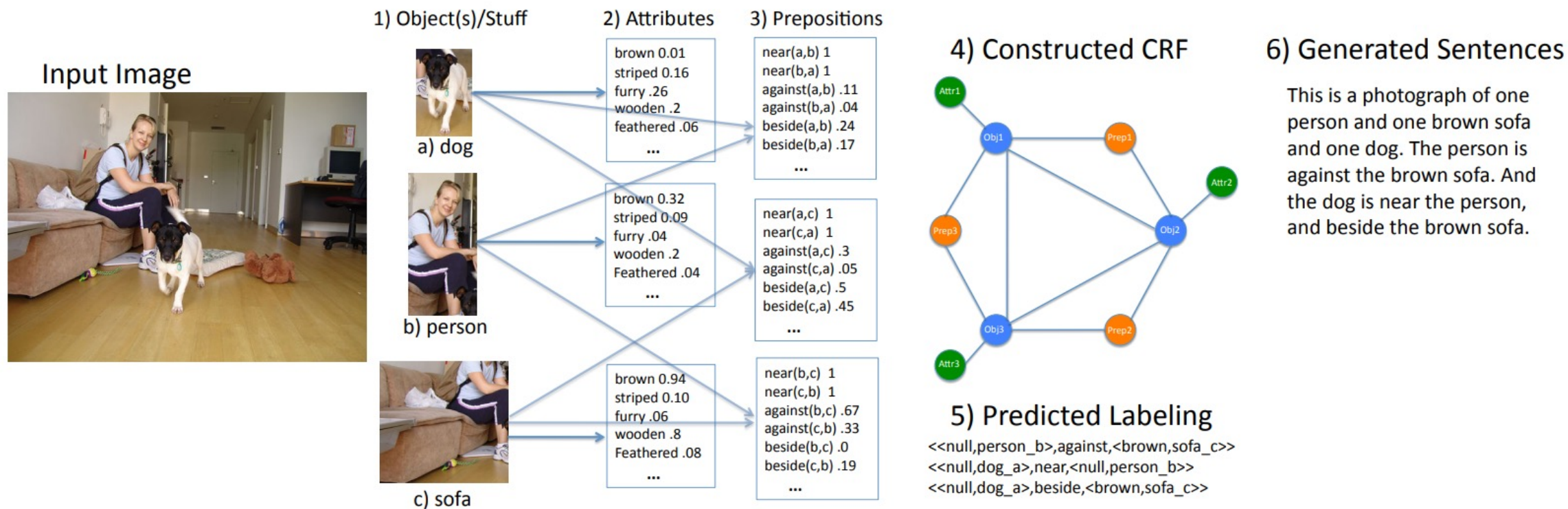


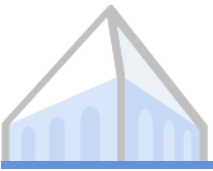
“Solved” Perception

- **Pro:** In early days of vision and language, assuming sub-problems provided traction.
- **Con:** Strong assumptions that don't hold in real world.

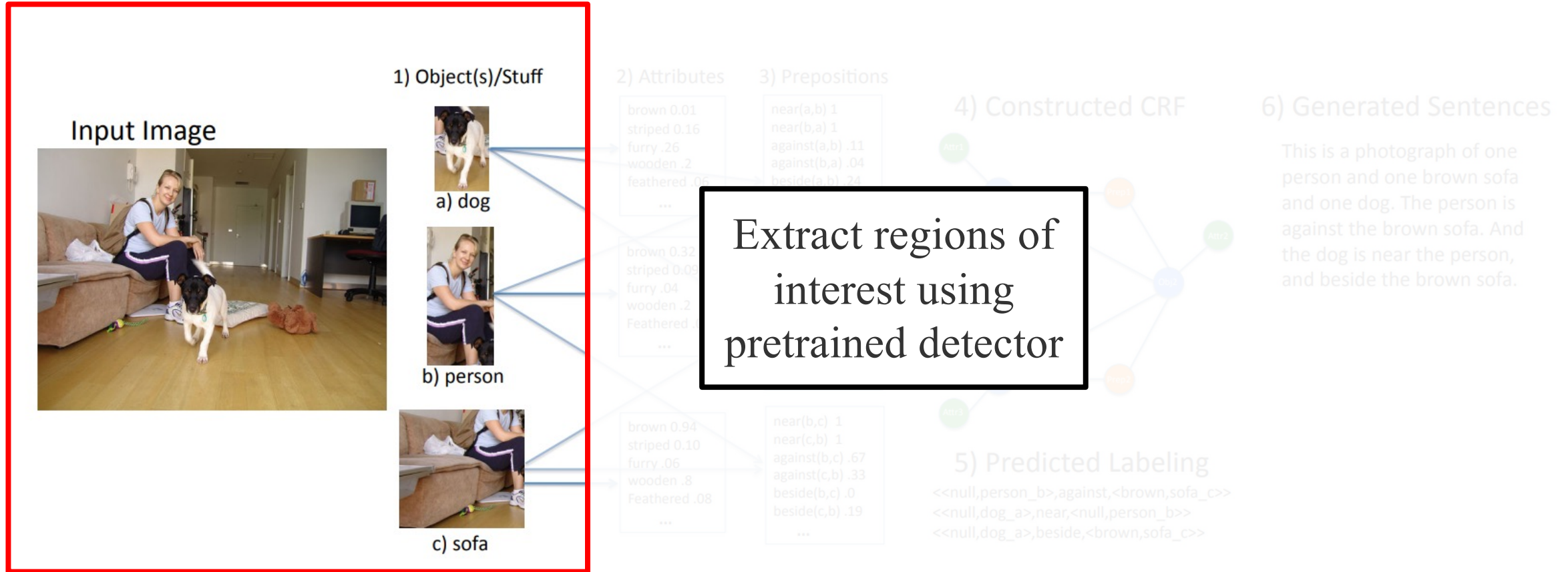


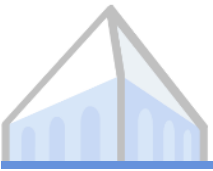
Intermediate Representations





Intermediate Representations





Intermediate Representations

Input Image



1) Object(s)/Stuff



a) dog



b) person



c) sofa

2) Attributes

brown 0.01
striped 0.16
furry .26
wooden .2
feathered .06
...

brown 0.32
striped 0.09
furry .04
wooden .2
Feathered .04
...

brown 0.94
striped 0.10
furry .06
wooden .8
Feathered .08
...

3) Prepositions

near(a,b) 1
near(b,a) 1
against(a,b) .11
against(b,a) .04
beside(a,b) .24
beside(b,a) .17
...

near(a,c) 1
near(c,a) 1
against(a,c) .3
against(c,a) .05
beside(a,c) .5
beside(c,a) .45
...

near(b,c) 1
near(c,b) 1
against(b,c) .67
against(c,b) .33
beside(b,c) .0
beside(c,b) .19
...

4) Constructed CRF



<<null, person_b>, against, <brown, sofa_c>>
<<null, dog_a>, near, <null, person_b>>
<<null, dog_a>, beside, <brown, sofa_c>>

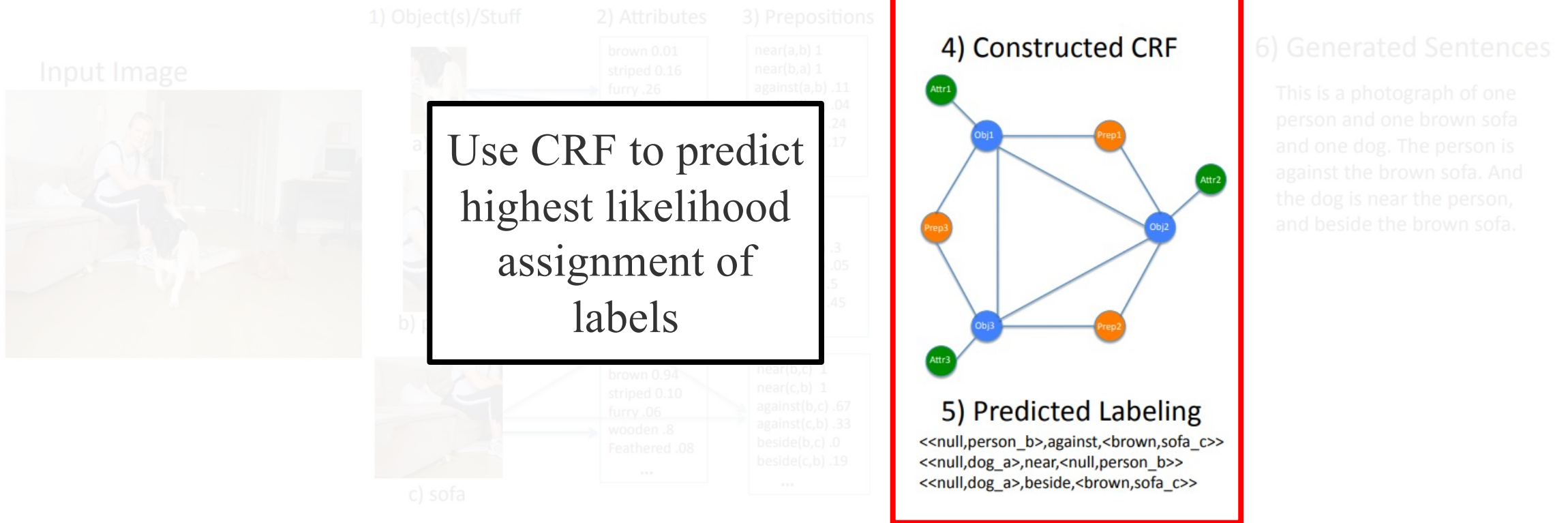
6) Generated Sentences

This is a photograph of one
the brown sofa
The person is
own sofa. And
for the person,
the brown sofa.

Classifiers score
attributes for each
region and
relationships
across them

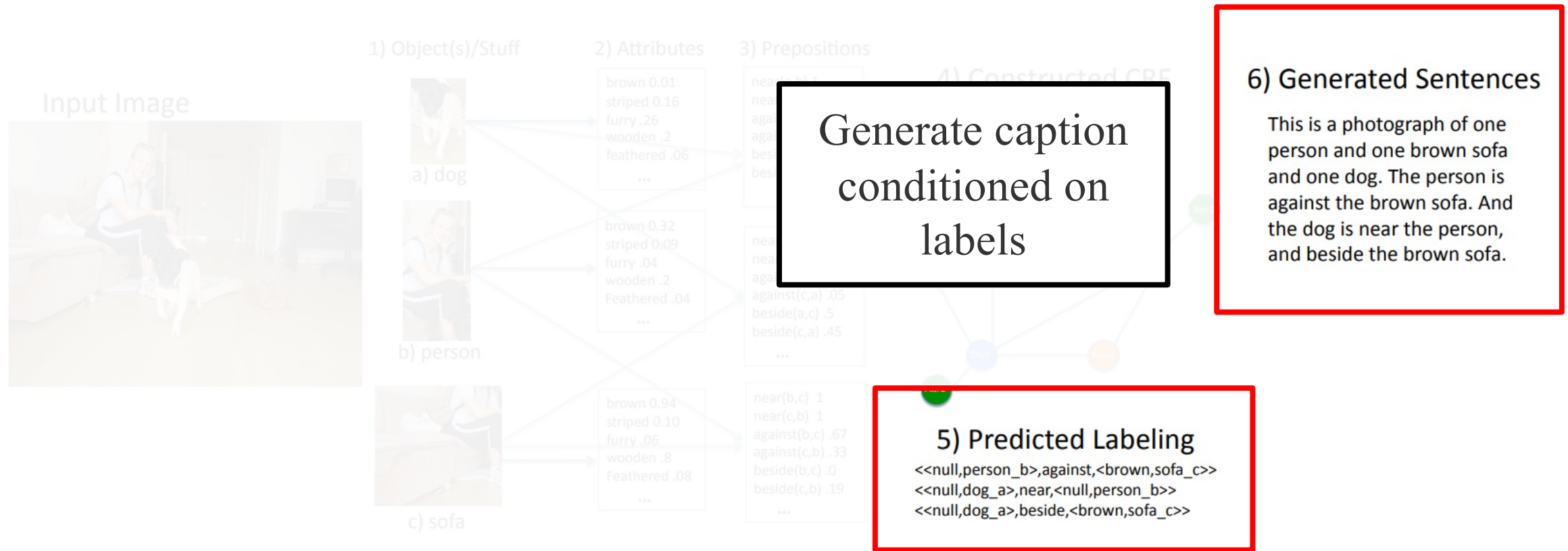


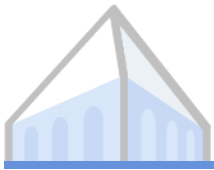
Intermediate Representations





Intermediate Representations

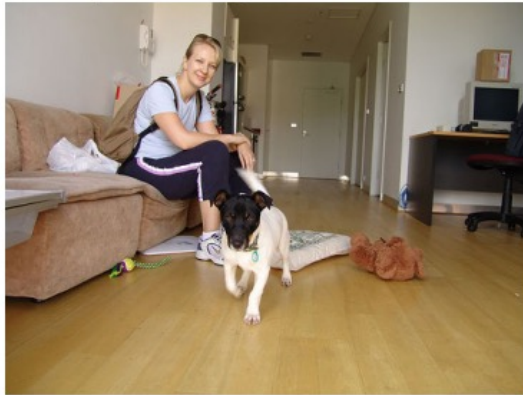




Intermediate Representations

Language model never sees pixels!

Input Image



1) Object(s)/Stuff



a) dog



b) person



c) sofa

2) Attributes

brown 0.01
striped 0.16
furry .26
wooden .2
feathered .06
...

brown 0.32
striped 0.09
furry .04
wooden .2
Feathered .04
...

brown 0.94
striped 0.10
furry .06
wooden .8
Feathered .08
...

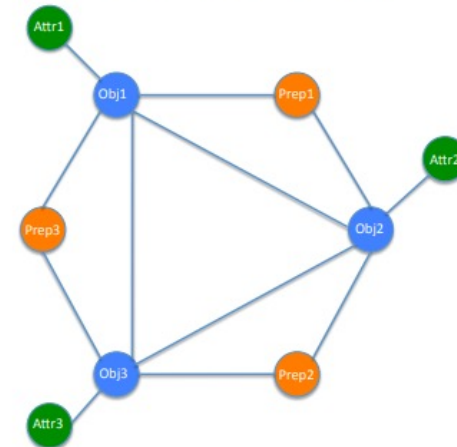
3) Prepositions

near(a,b) 1
near(b,a) 1
against(a,b) .11
against(b,a) .04
beside(a,b) .24
beside(b,a) .17
...

near(a,c) 1
near(c,a) 1
against(a,c) .3
against(c,a) .05
beside(a,c) .5
beside(c,a) .45
...

near(b,c) 1
near(c,b) 1
against(b,c) .67
against(c,b) .33
beside(b,c) .0
beside(c,b) .19
...

4) Constructed CRF

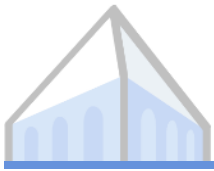


5) Predicted Labeling

<<null, person_b>, against, <brown, sofa_c>>
<<null, dog_a>, near, <null, person_b>>
<<null, dog_a>, beside, <brown, sofa_c>>

6) Generated Sentences

This is a photograph of one person and one brown sofa and one dog. The person is against the brown sofa. And the dog is near the person, and beside the brown sofa.



Intermediate Representations



This is a photograph of one sky, one road and one bus. The blue sky is above the gray road. The gray road is near the shiny bus. The shiny bus is near the blue sky.



There are two aeroplanes. The first shiny aeroplane is near the second shiny aeroplane.



There are one cow and one sky. The golden cow is by the blue sky.



There are one dining table, one chair and two windows. The wooden dining table is by the wooden chair, and against the first window, and against the second white window. The wooden chair is by the first window, and by the second white window. The first window is by the second white window.



Here we see one person and one train. The black person is by the train.



This is a picture of one sky, one road and one sheep. The gray sky is over the gray road. The gray sheep is by the gray road.



Here we see one road, one sky and one bicycle. The road is near the blue sky, and near the colorful bicycle. The colorful bicycle is within the blue sky.



Here we see two persons, one sky and one aeroplane. The first black person is by the blue sky. The blue sky is near the shiny aeroplane. The second black person is by the blue sky. The shiny aeroplane is by the first black person, and by the second black person.



This is a picture of two dogs. The first dog is near the second furry dog.



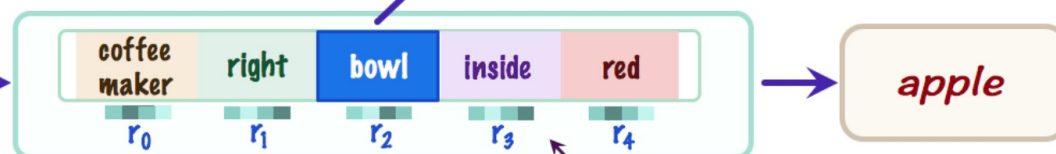
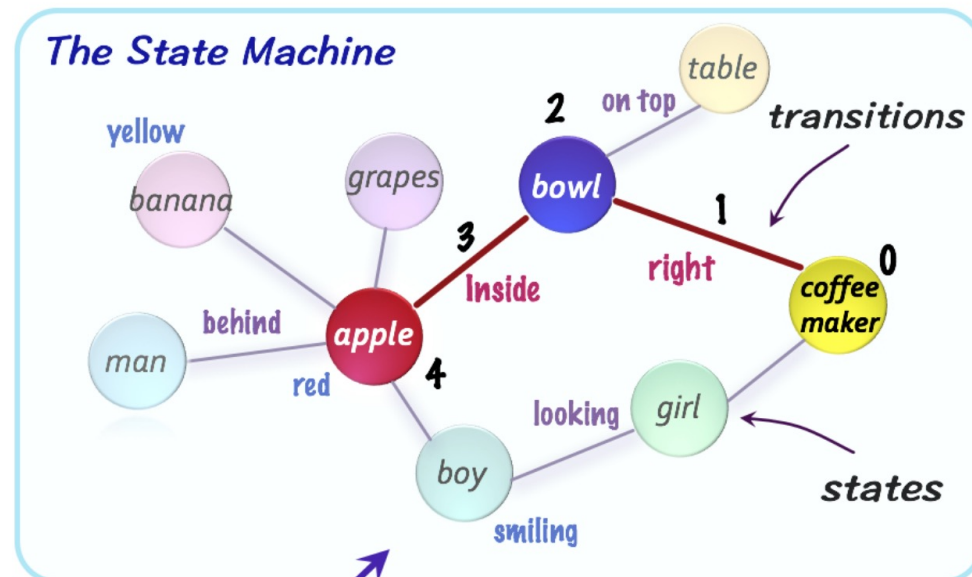
This is a photograph of two buses. The first rectangular bus is near the second rectangular bus.



Intermediate Representations



What is the **red** fruit **inside** the **bowl** to the **right** of the **coffee maker**?



instructions

apple

alphabet (concepts)

bowl

Color: brown (0.92)

Material: wood (0.8)

apple

Color: red (0.95)

Shape: round (0.87)

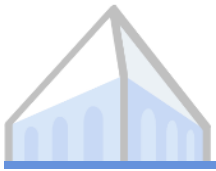
girl

Mood: happy (0.78)

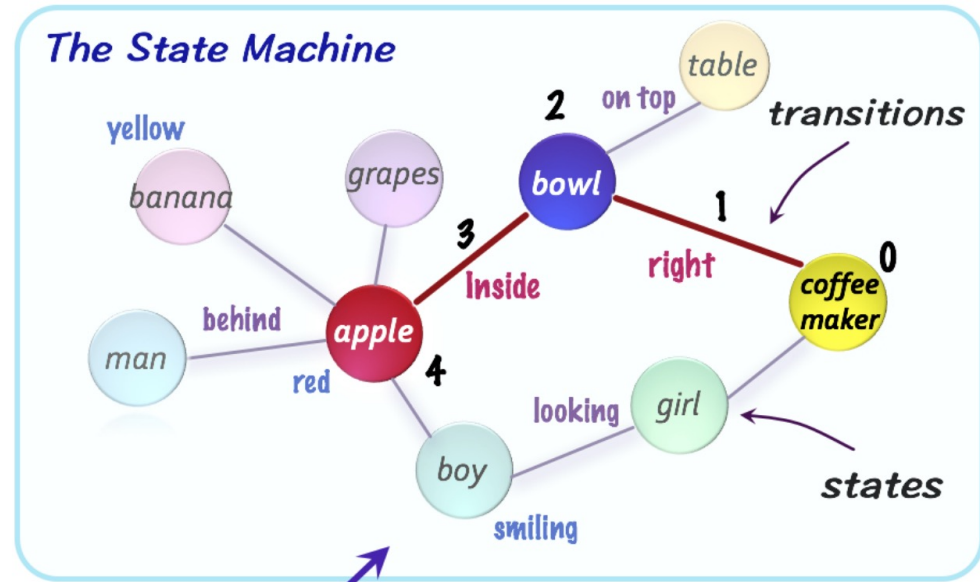
Posture: sitting (0.82)

properties

disentangled
representation



Intermediate Representations



alphabet (concepts)

Generate scene graph from image

What is the red fruit inside the bowl to the right of the coffee maker?



instructions

properties

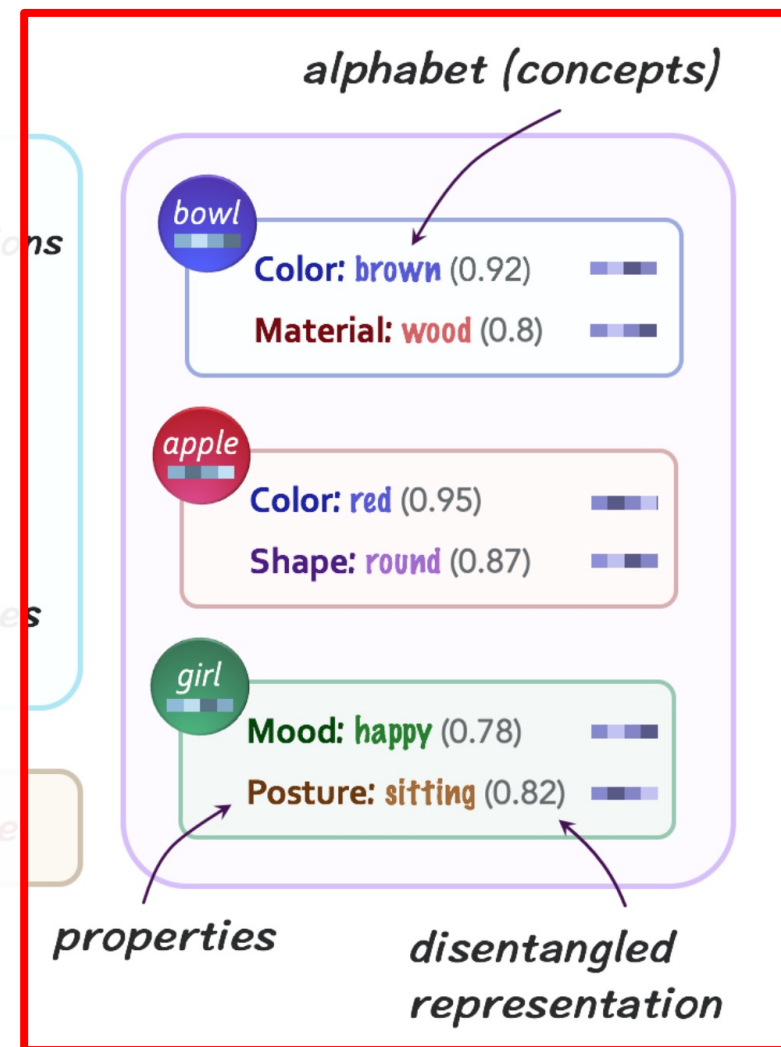
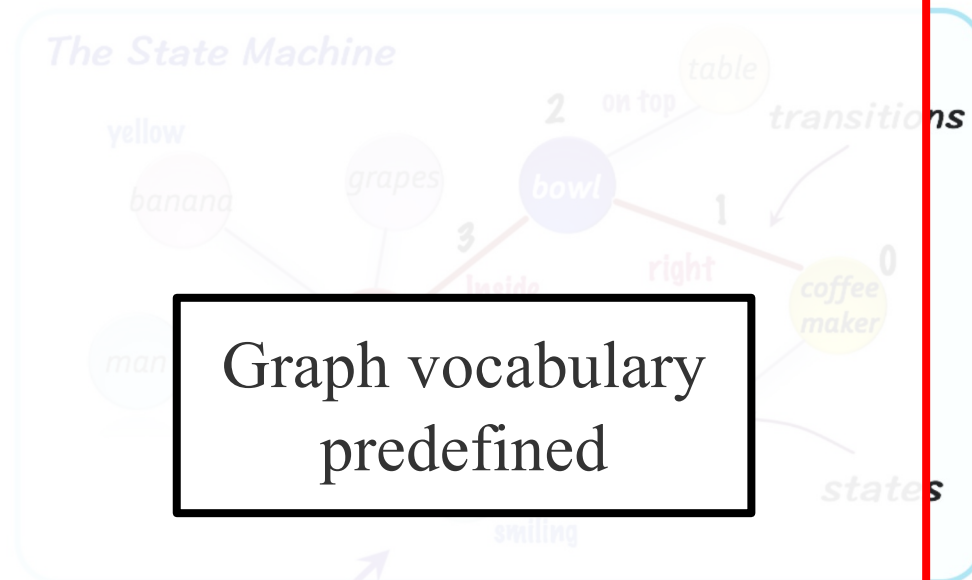
disentangled representation

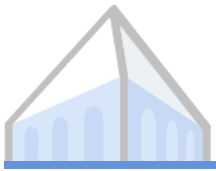


Intermediate Representations



What is the red fruit inside the bowl to the right of the coffee maker?

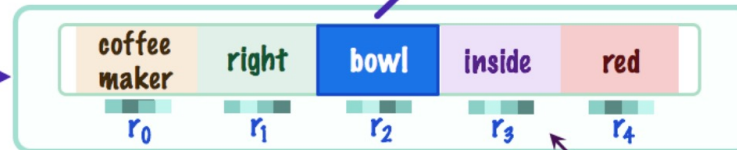




Intermediate Representations

Transform
question into
program traversing
graph for answer

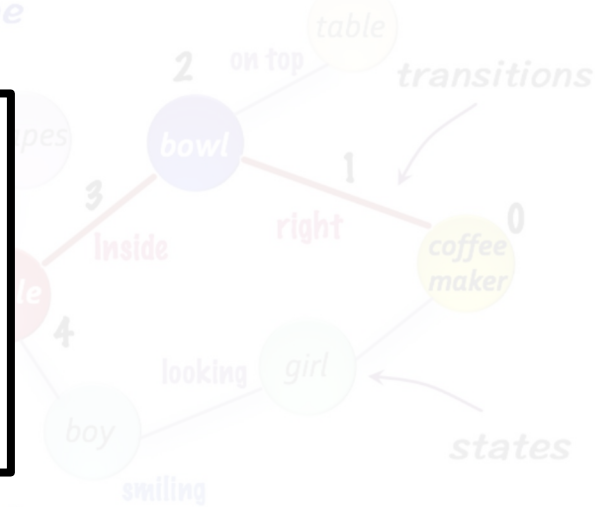
What is the *red fruit* inside the *bowl*
to the *right* of the *coffee maker*?



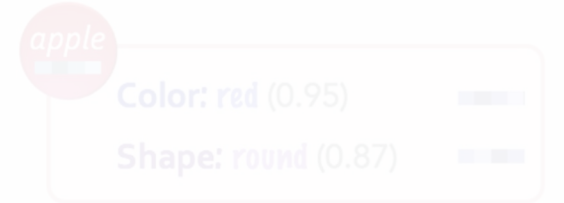
instructions

apple

The State Machine



alphabet (concepts)



properties

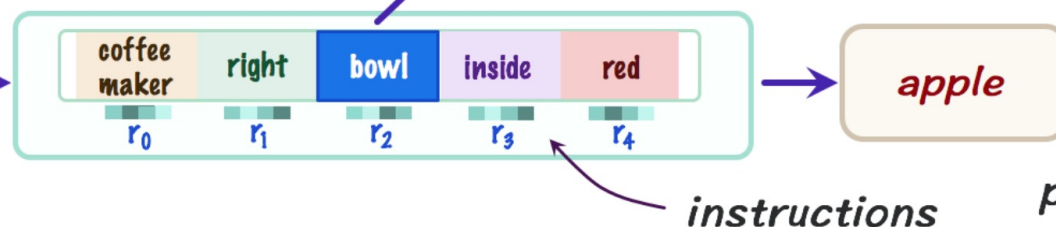
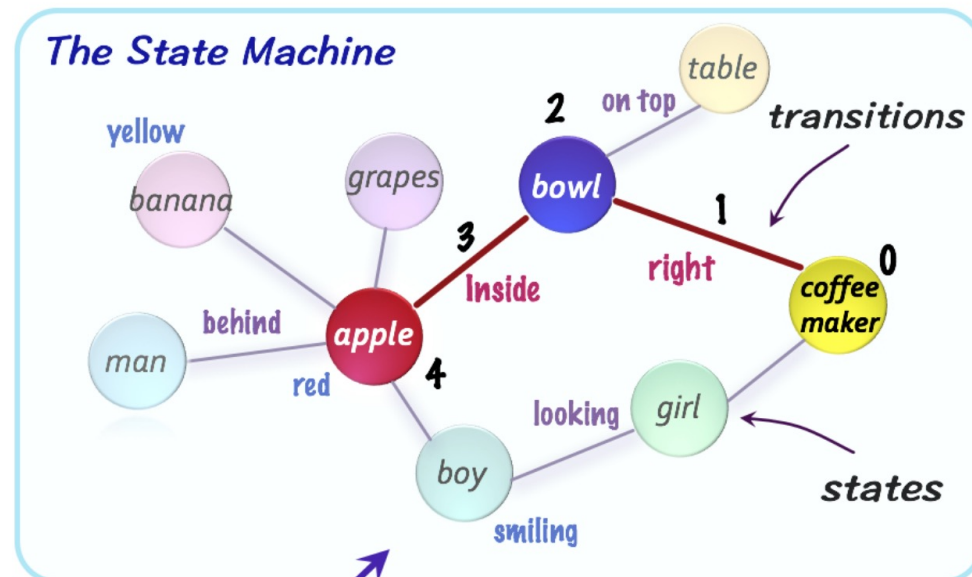
disentangled
representation



Intermediate Representations



What is the *red fruit* inside the *bowl* to the *right* of the *coffee maker*?



alphabet (concepts)



Answer by
executing program
in state machine

Mood: happy (0.78)
Posture: sitting (0.82)

properties

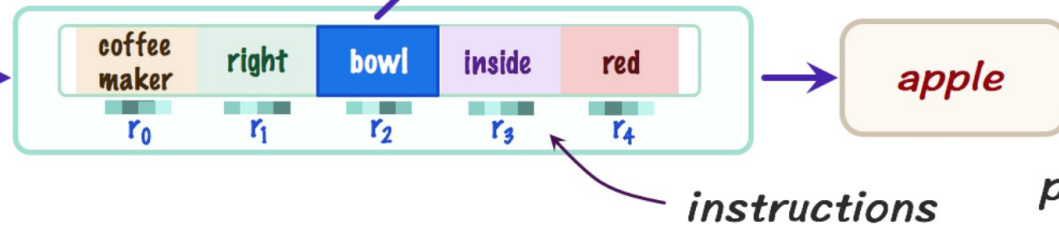
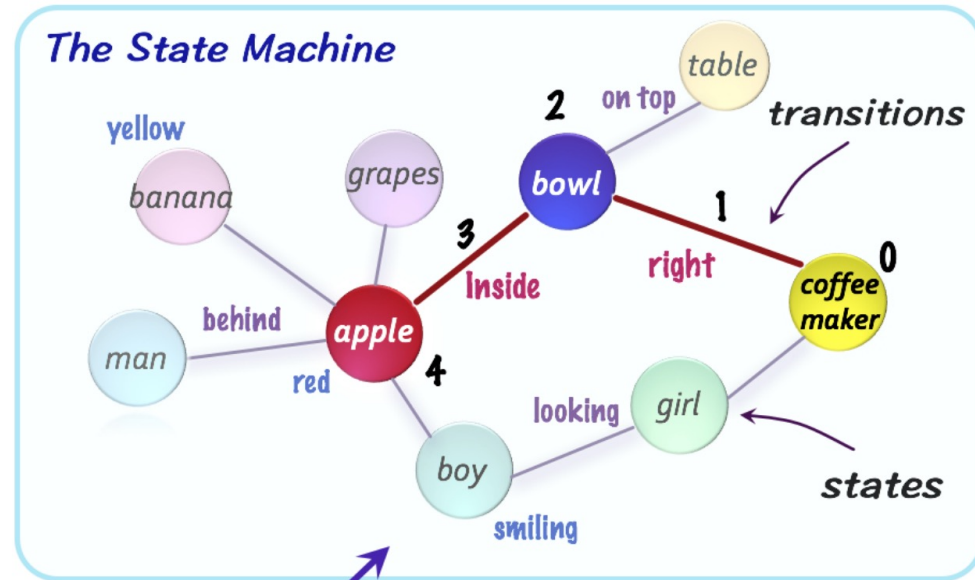
*disentangled
representation*



Intermediate Representations



What is the *red fruit* inside the *bowl* to the *right* of the *coffee maker*?



alphabet (concepts)

bowl
Color: brown (0.92)
Material: wood (0.8)

Allows language reasoning to occur solely within abstract structure

Posture: sitting (0.82)

properties

disentangled representation



Intermediate Representations

Table 4: GQA generalization

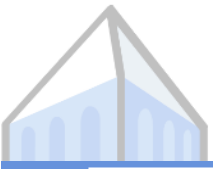
Model	Content	Structure
Global Prior	8.51	14.64
Local Prior	12.14	18.21
Vision	17.51	18.68
Language	21.14	32.88
Lang+Vis	24.95	36.51
BottomUp [5]	29.72	41.83
MAC [40]	31.12	47.27
NSM	40.24	55.72

TEXT PROMPT

a store front that has the word 'dall-e' written on it. a store front that has the word 'dall-e' written on it. a store front that has the word 'dall-e' written on it. dall-e store front.

AI-GENERATED
IMAGES





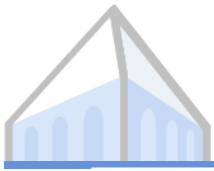
Intermediate Representations

Step 1

Learn Proto-linguistic
Code Book



1	5	2	6
9	13	10	14
3	7	4	8
11	15	12	16



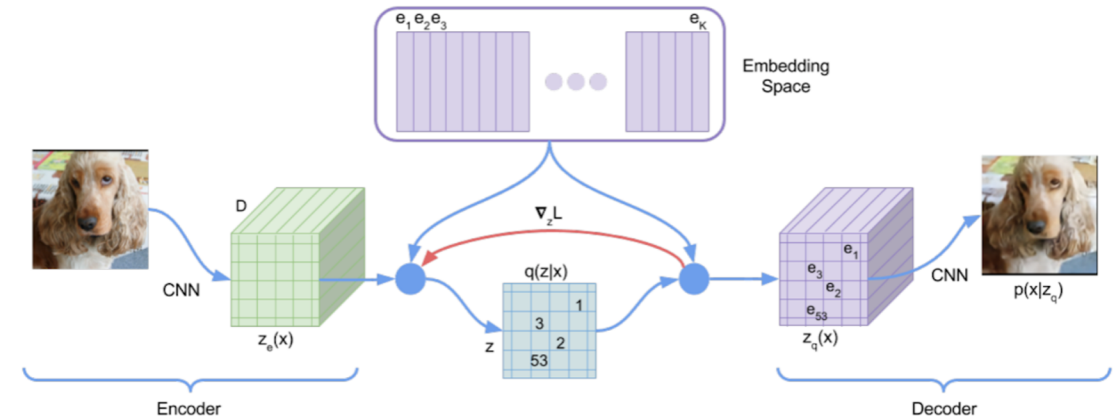
Intermediate Representations

Step 1

Learn Proto-linguistic
Code Book



1	5	2	6
9	13	10	14
3	7	4	8
11	15	12	16



Neural Discrete Representation Learning: van Oord et al. 2017



Intermediate Representations

Step 2

Learn Joint

Language and Code Distribution

"A kitten
with a pink
background"

1	5	2	6
9	13	10	14
3	7	4	8
11	15	12	16



Intermediate Representations

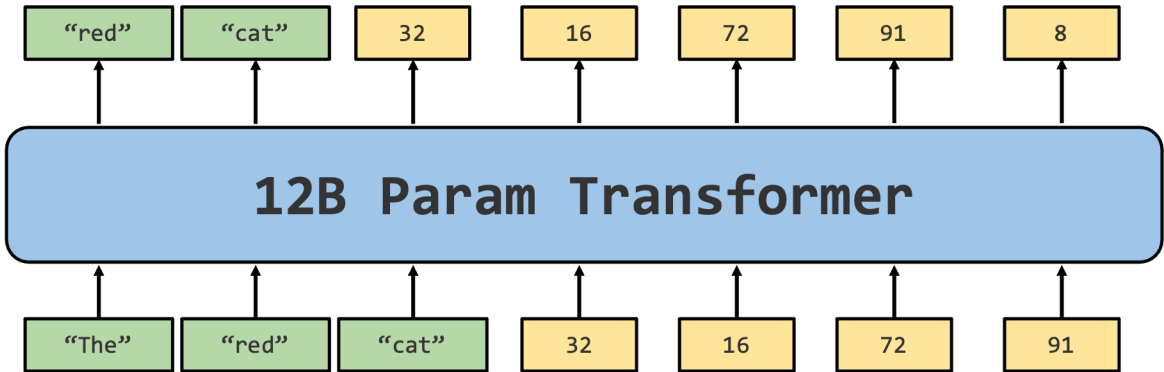
Step 2

Learn Joint

Language and Code Distribution

"A kitten
with a pink
background"

1	5	2	6
9	13	10	14
3	7	4	8
11	15	12	16



Generating Long Sequences with Sparse Transformers: Child et al. 2019



Intermediate Representations

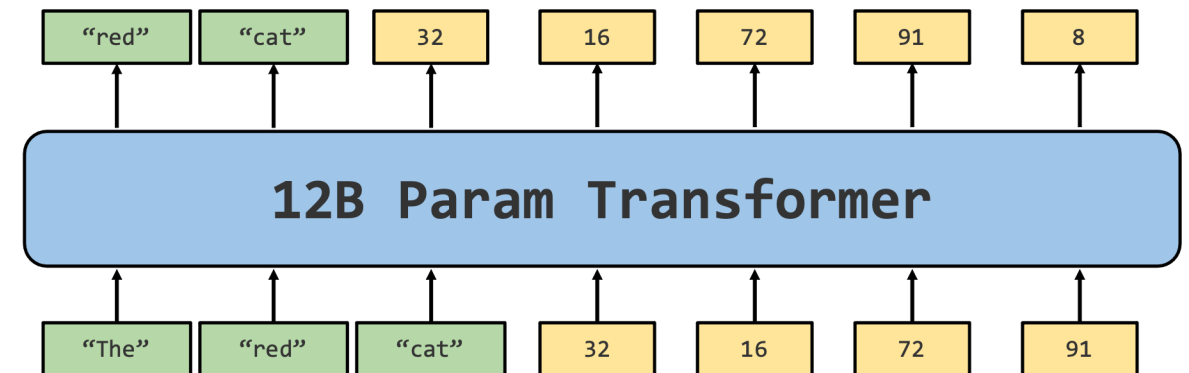
Step 2

Learn Joint

Language and Code Distribution

"A kitten
with a pink
background"

1	5	2	6
9	13	10	14
3	7	4	8
11	15	12	16



Generating Long Sequences with Sparse Transformers: Child et al. 2019

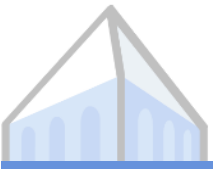
Reduced to language modeling
problem!

TEXT PROMPT

an x-ray of a capybara sitting in a forest

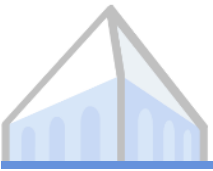
AI-GENERATED
IMAGES





Anchoring to 3D

*“The goal of an image understanding system is to transform two-dimensional data into a **representation** of the three-dimensional spatio-temporal world”*



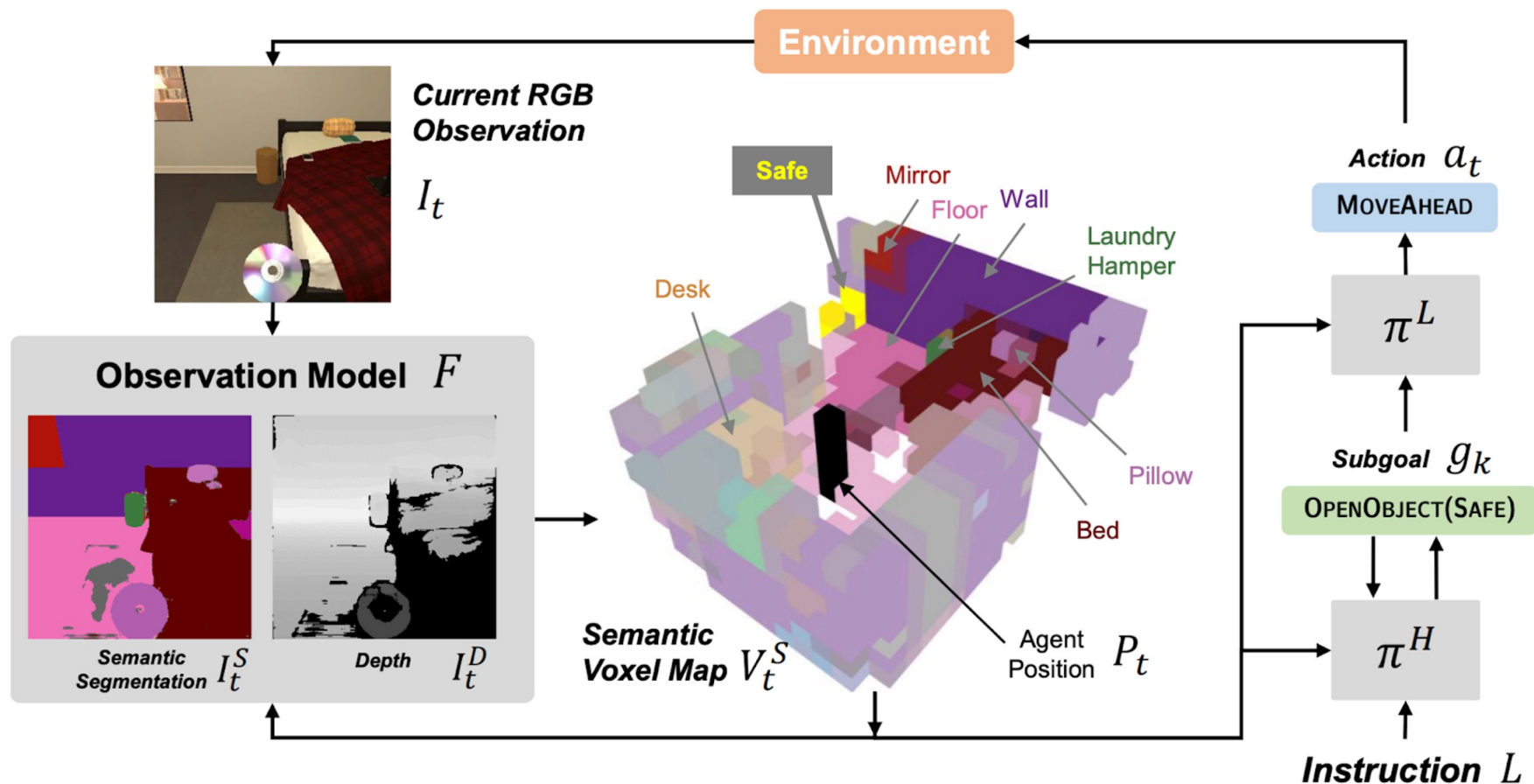
Anchoring to 3D



“Place a clean ladle on a counter”

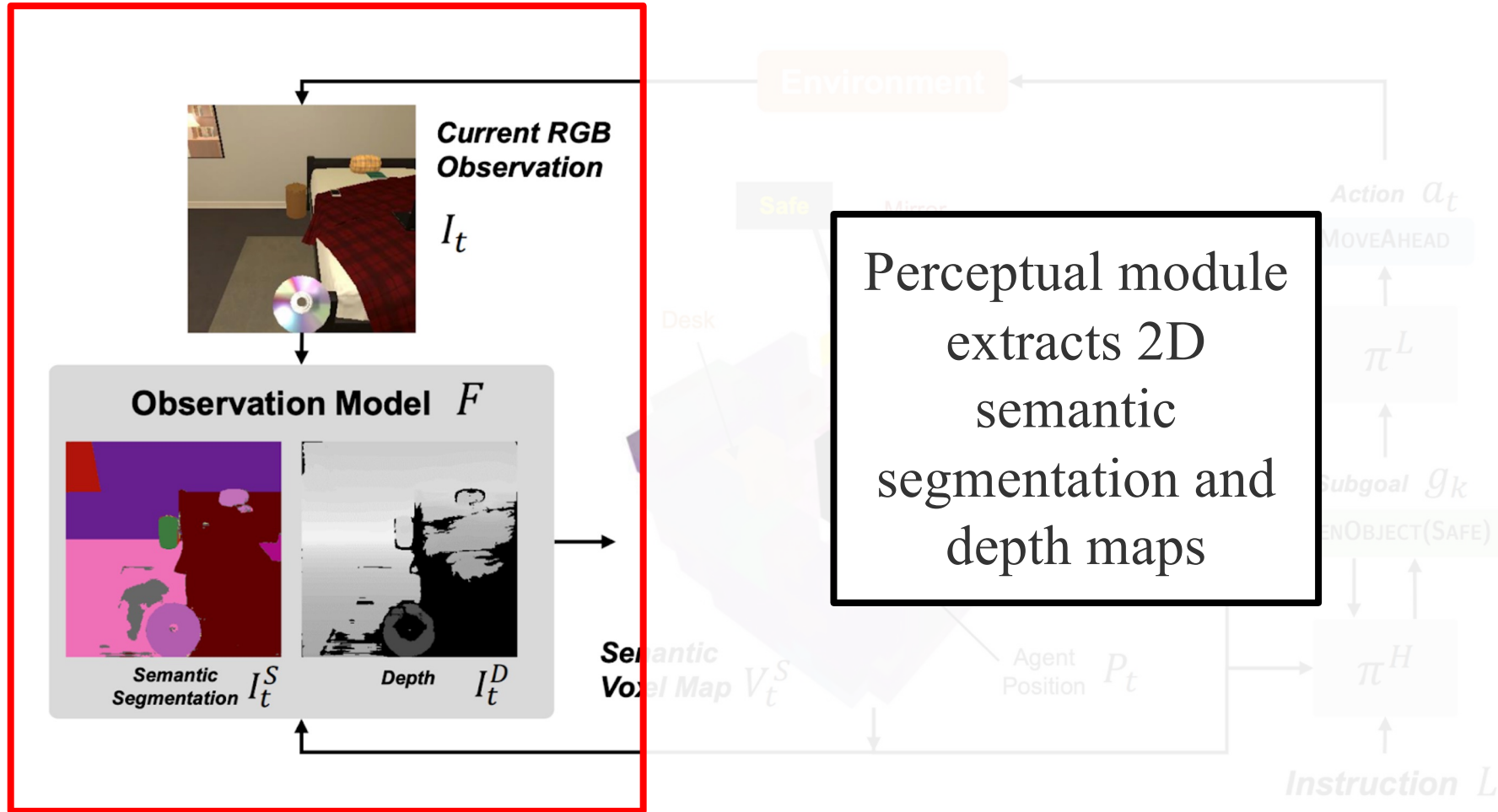


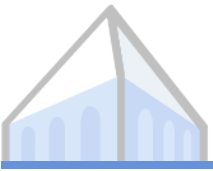
Anchoring to 3D



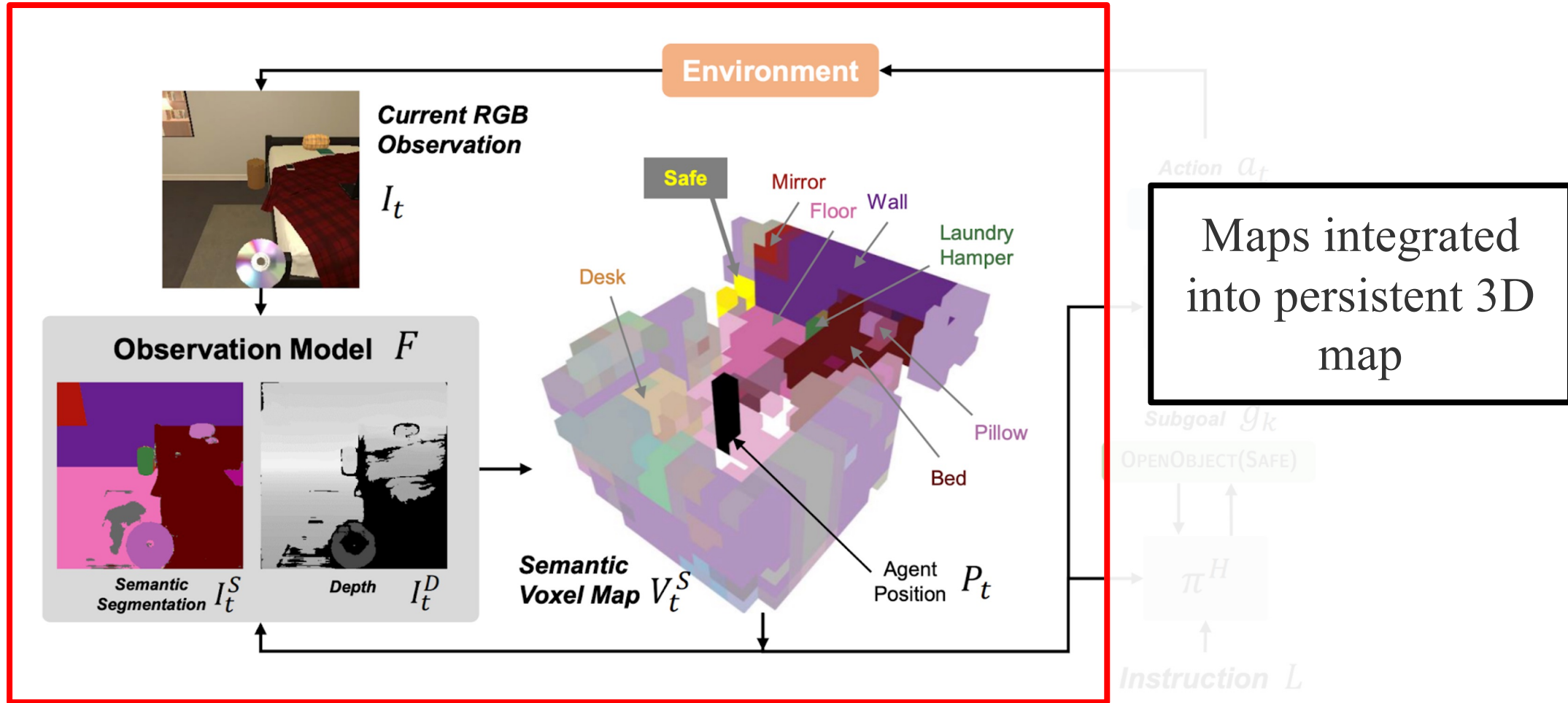


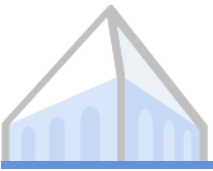
Anchoring to 3D



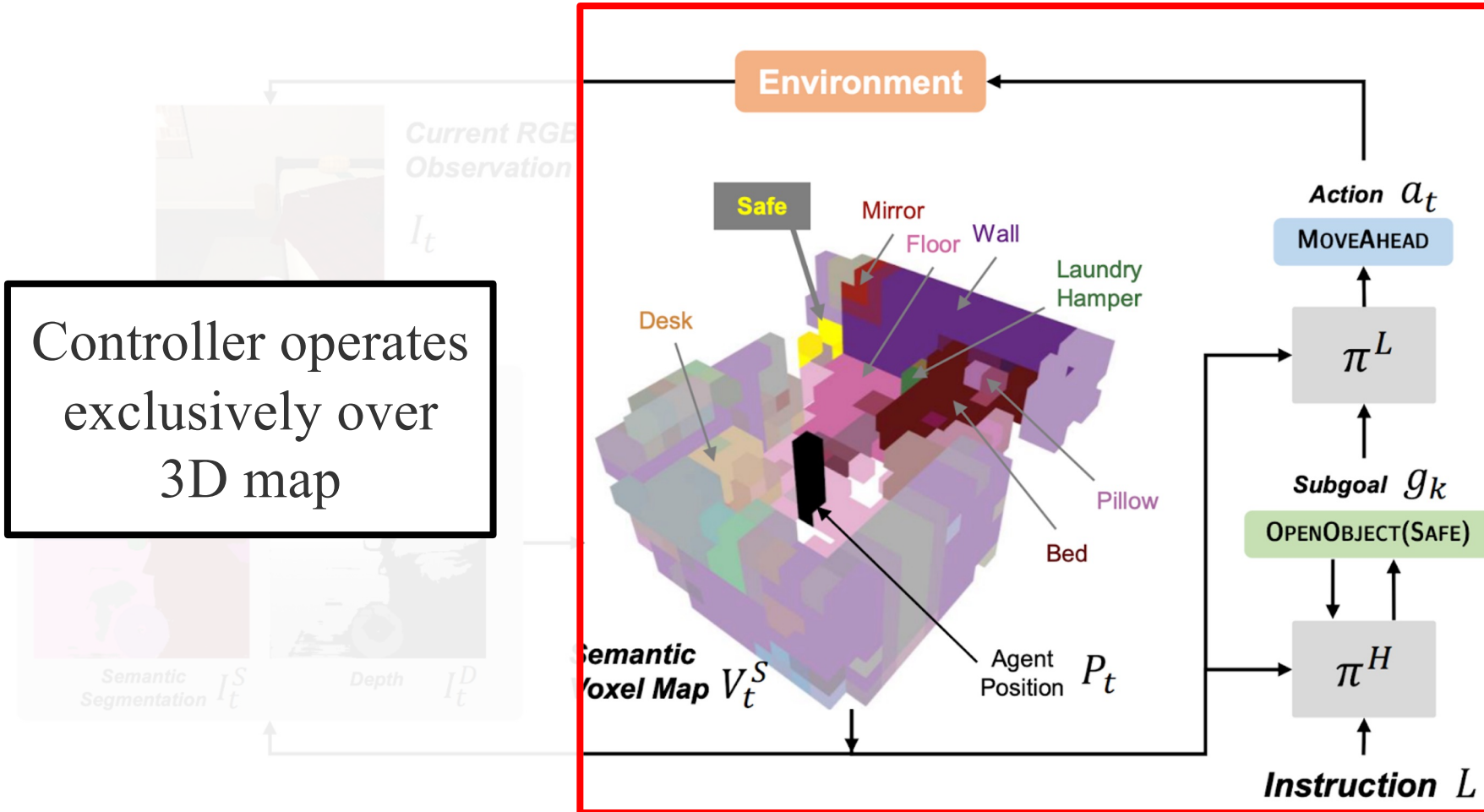


Anchoring to 3D





Anchoring to 3D





Anchoring to 3D

Method	Validation			
	Seen		Unseen	
	SR	GC	SR	GC
HLSM	29.6	38.8	18.3	31.2
+ gt depth	29.6	40.5	20.1	33.7
+ gt depth, gt seg.	40.7	50.4	40.2	52.2
+ gt seg.	36.2	47.0	34.7	47.8
w/o language enc.	0.9	8.6	0.2	7.5
w/o subg. hist. enc.	29.4	38.5	16.6	29.2
w/o state repr enc.	30.0	40.6	18.9	30.8

3D Map
useful for
improving
performance



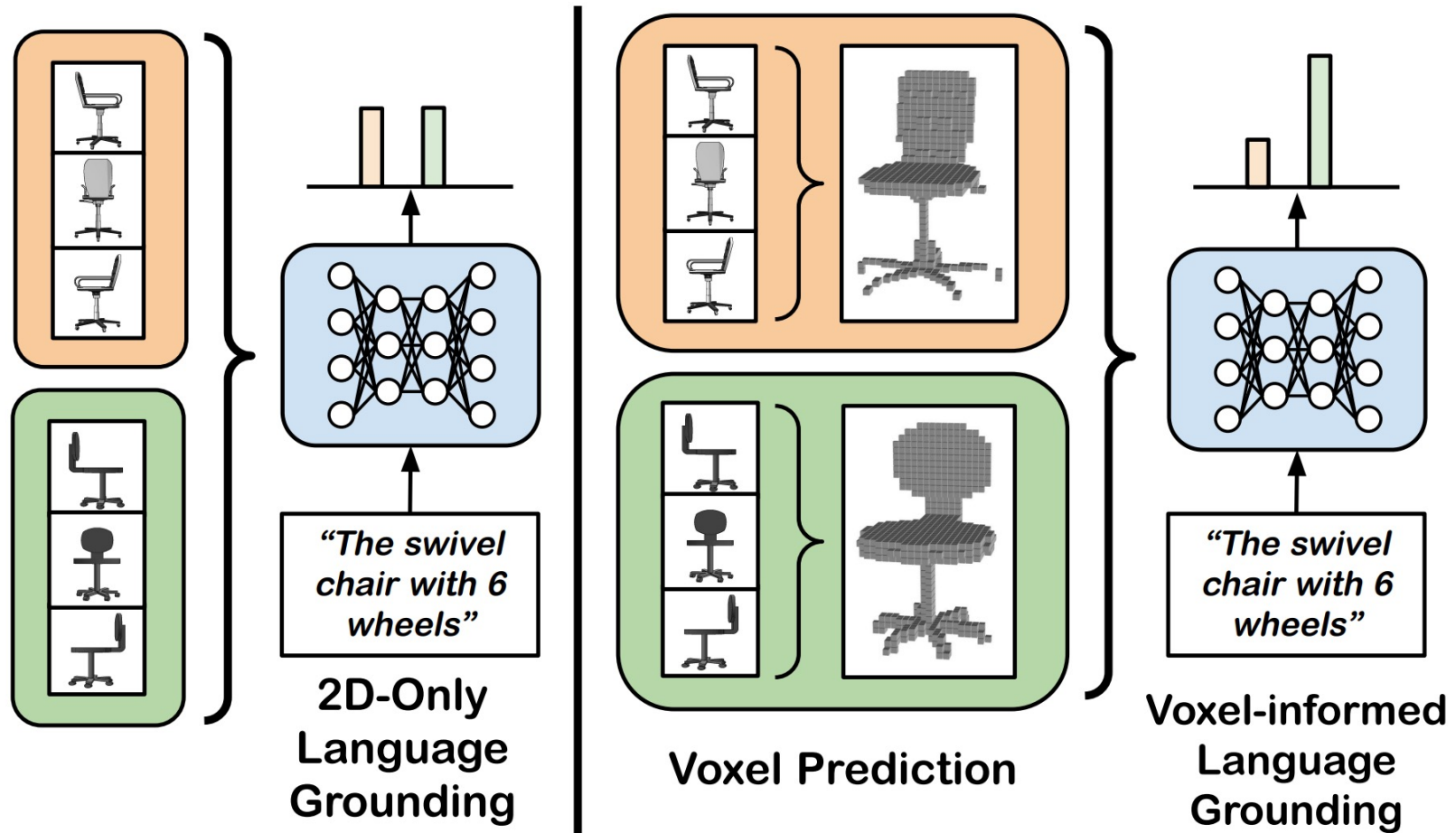
Anchoring to 3D

Method	Validation			
	Seen		Unseen	
	SR	GC	SR	GC
HLSM	29.6	38.8	18.3	31.2
+ gt depth	29.6	40.5	20.1	33.7
+ gt depth, gt seg.	40.7	50.4	40.2	52.2
+ gt seg.	36.2	47.0	34.7	47.8
w/o language enc.	0.9	8.6	0.2	7.5
w/o subg. hist. enc.	29.4	38.5	16.6	29.2
w/o state repr enc.	30.0	40.6	18.9	30.8

However,
benefits held
back by
cascading
errors

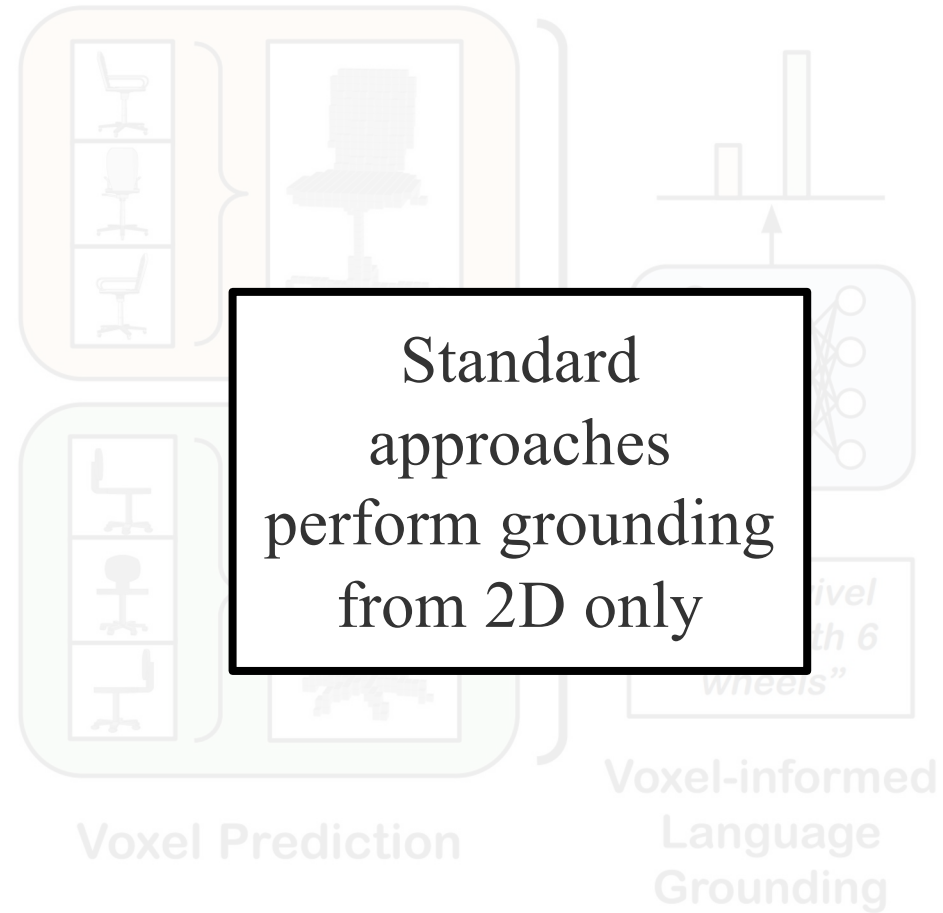
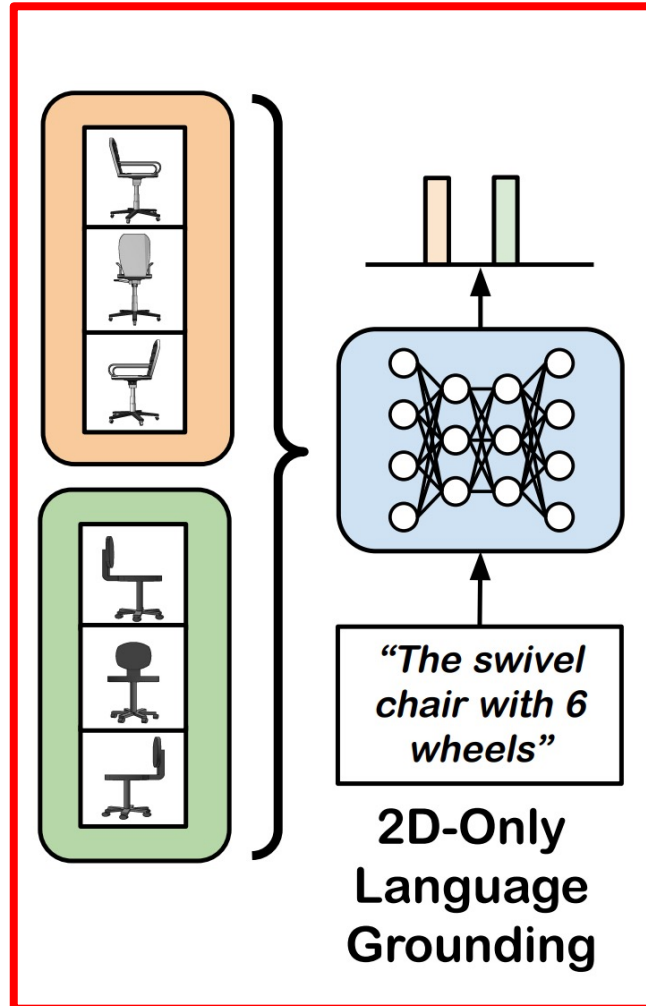


Anchoring to 3D

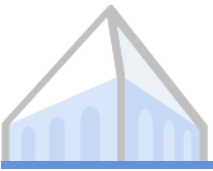




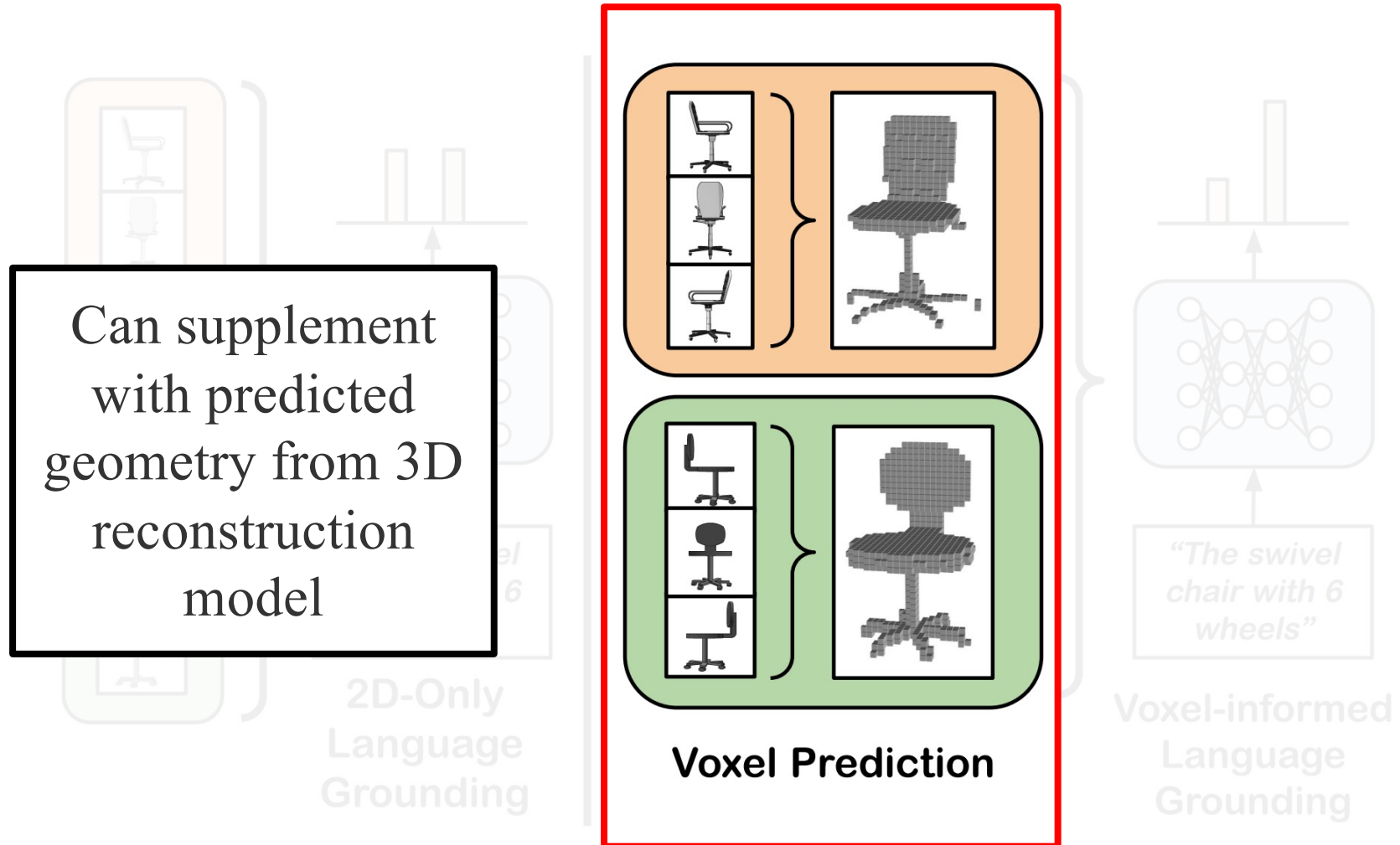
Anchoring to 3D



Standard
approaches
perform grounding
from 2D only



Anchoring to 3D





Anchoring to 3D

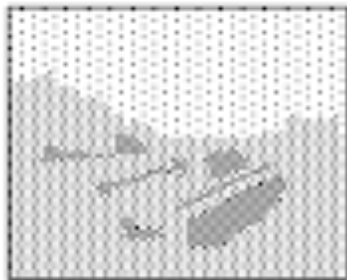
Model	VALIDATION			TEST		
	Visual	Blind	All	Visual	Blind	All
ViLBERT	89.5	76.6	83.1	80.2	73.0	76.6
MATCH	89.2 (0.9)	75.2 (0.7)	82.2 (0.4)	83.9 (0.5)	68.7 (0.9)	76.5 (0.5)
MATCH*	90.6 (0.4)	75.7 (1.2)	83.2 (0.8)	-	-	-
LAGOR	89.8 (0.4)	75.3 (0.7)	82.6 (0.4)	84.3 (0.4)	69.4 (0.5)	77.0 (0.5)
LAGOR*	89.8 (0.5)	75.0 (0.4)	82.5 (0.1)	-	-	-
VLG (Ours)	91.2 (0.4)	78.4[†] (0.7)	84.9[†] (0.3)	86.0	71.7	79.0

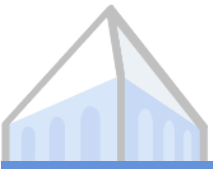
Improves performance over 2D-only
approaches



Bottom-Up Takeaways

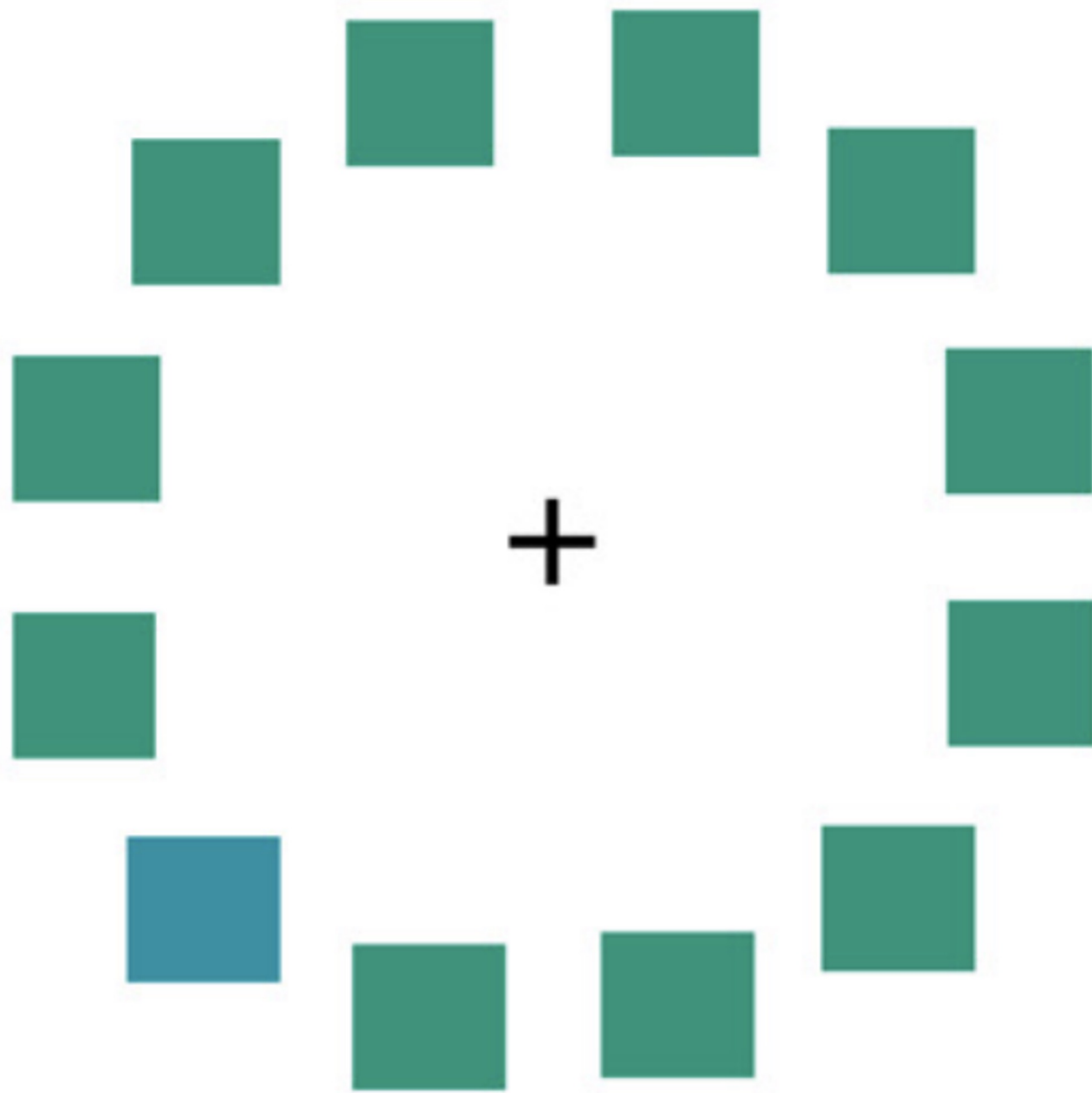
- Grounding to intermediate representations more tractable than grounding directly to pixels.
- Constrains the space of things to ground to.
- **Limitation:**
 - May suffer from cascading error.
 - Not always informed by language.

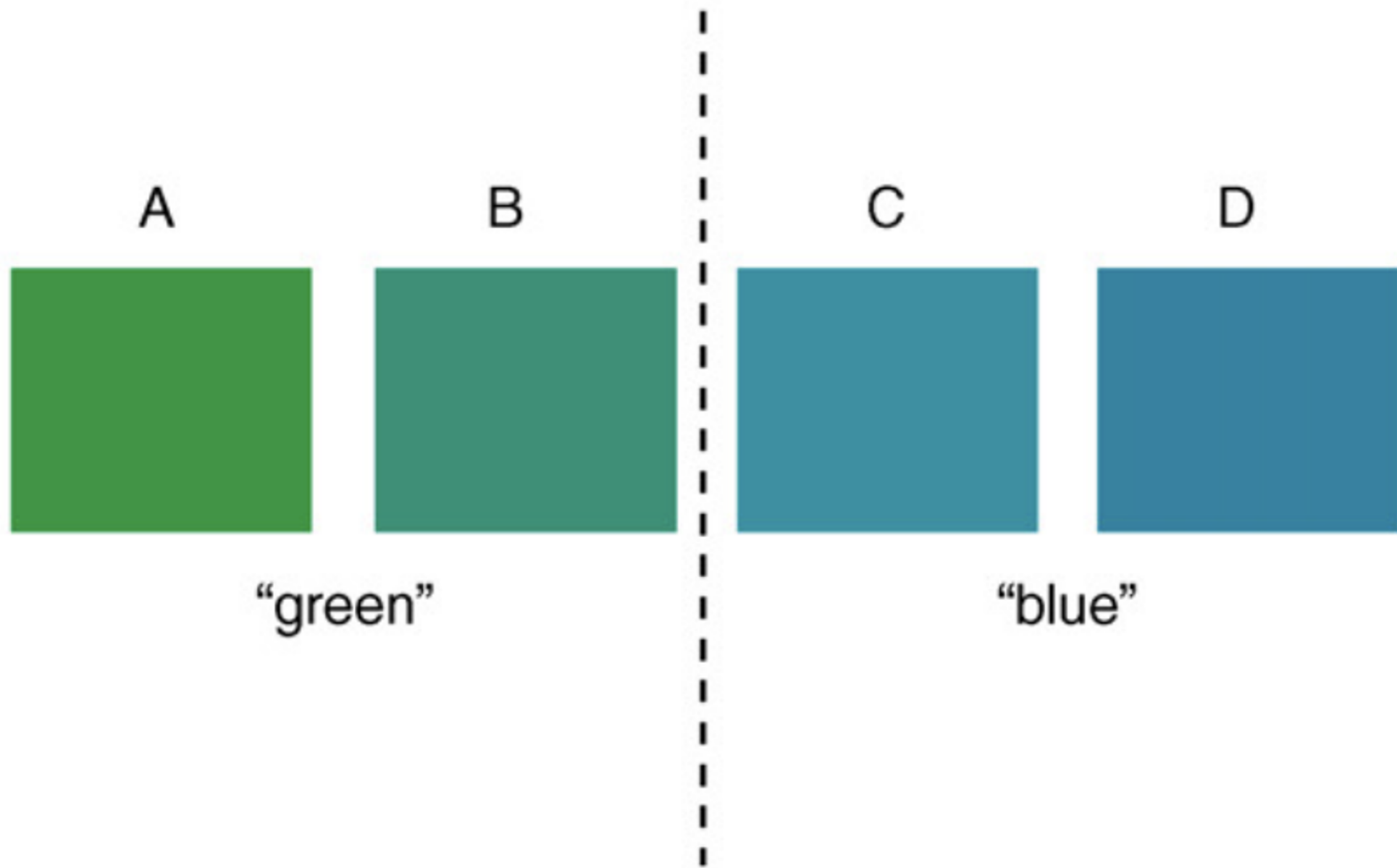




Top-Down

“What color is the
small **shiny cube**?”







WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

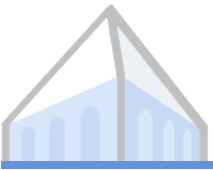
Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

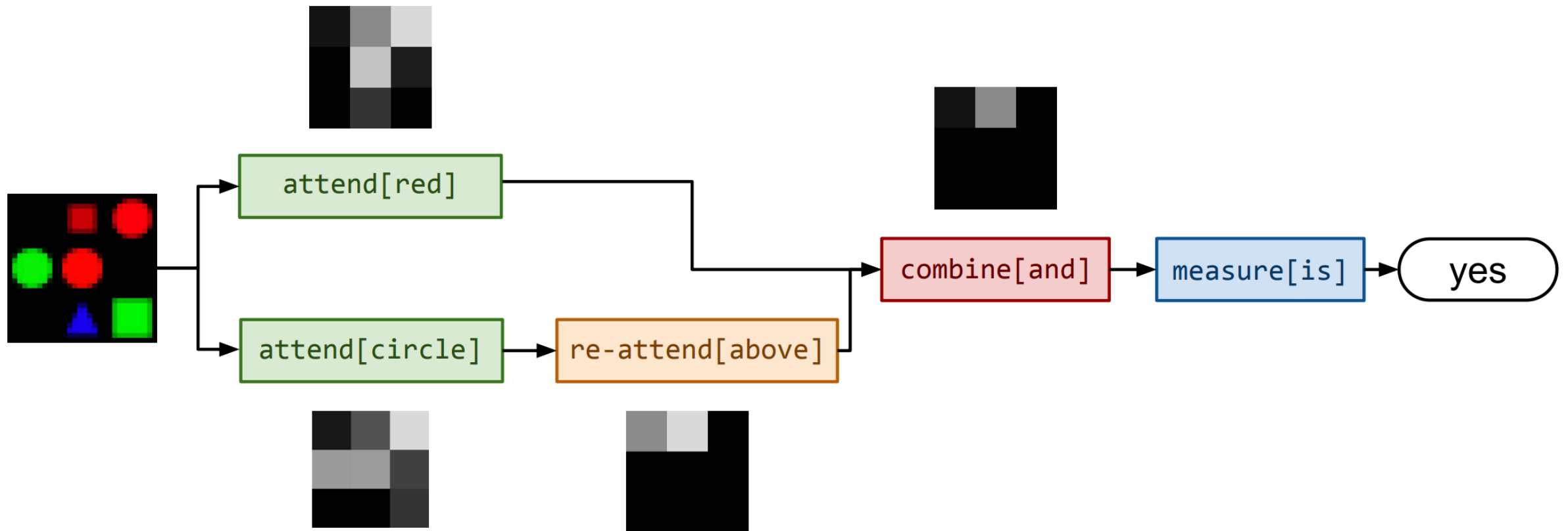
Noun

- [S:](#) (n) **wordnet** (any of the machine-readable lexical databases modeled after the Princeton WordNet)
- [S:](#) (n) **WordNet**, [Princeton WordNet](#) (a machine-readable lexical database organized by meanings; developed at Princeton University)



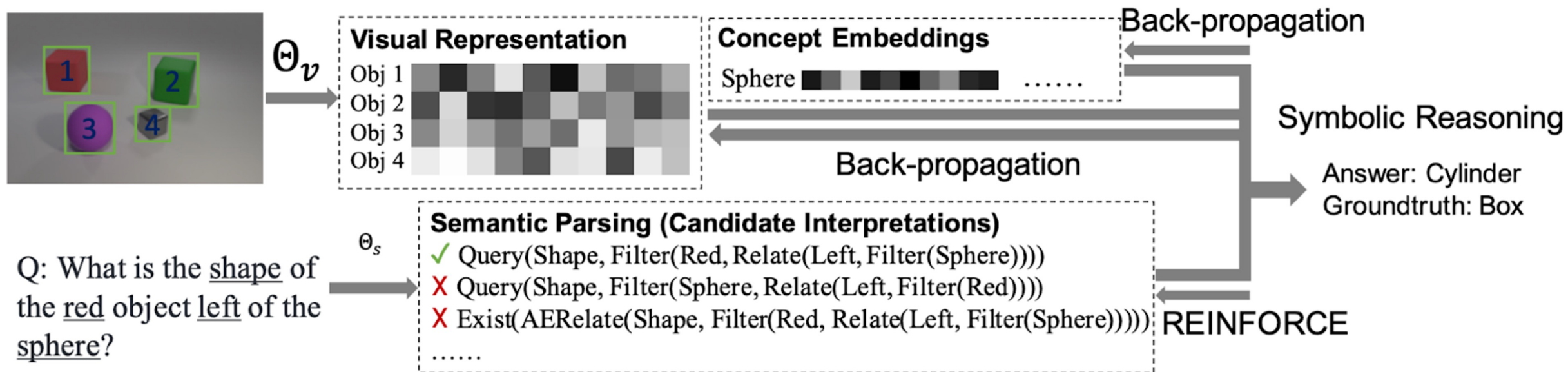
Modular Systems

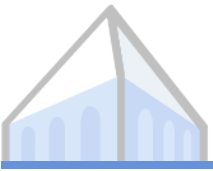
“Is there a red sphere above a circle?”



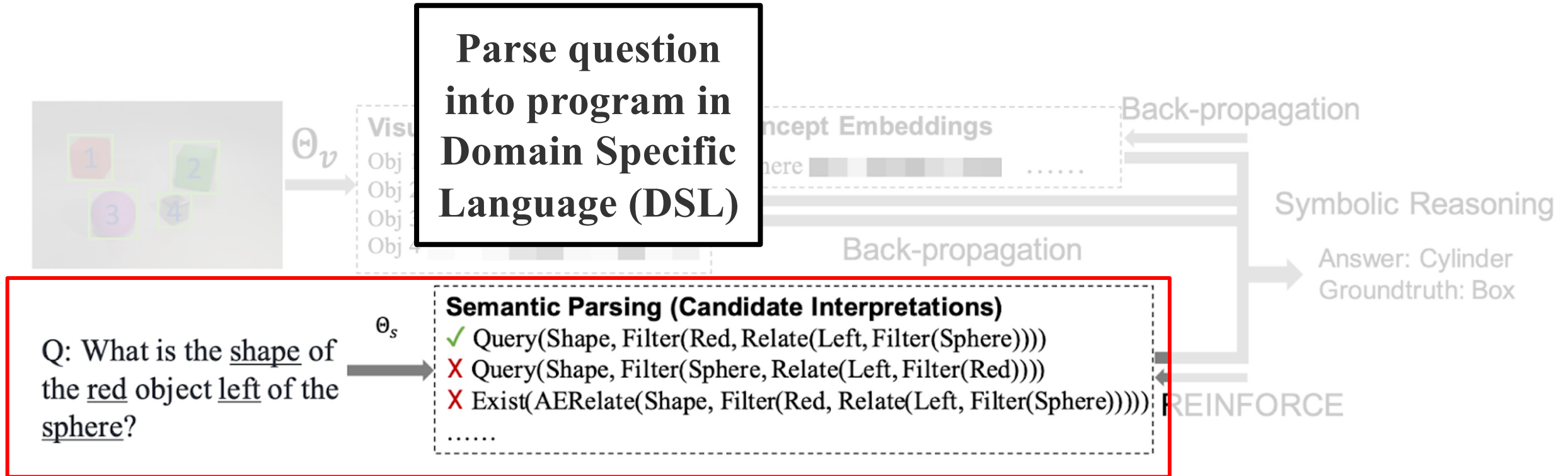


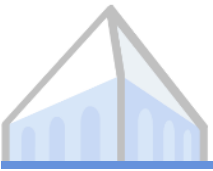
Modular Systems



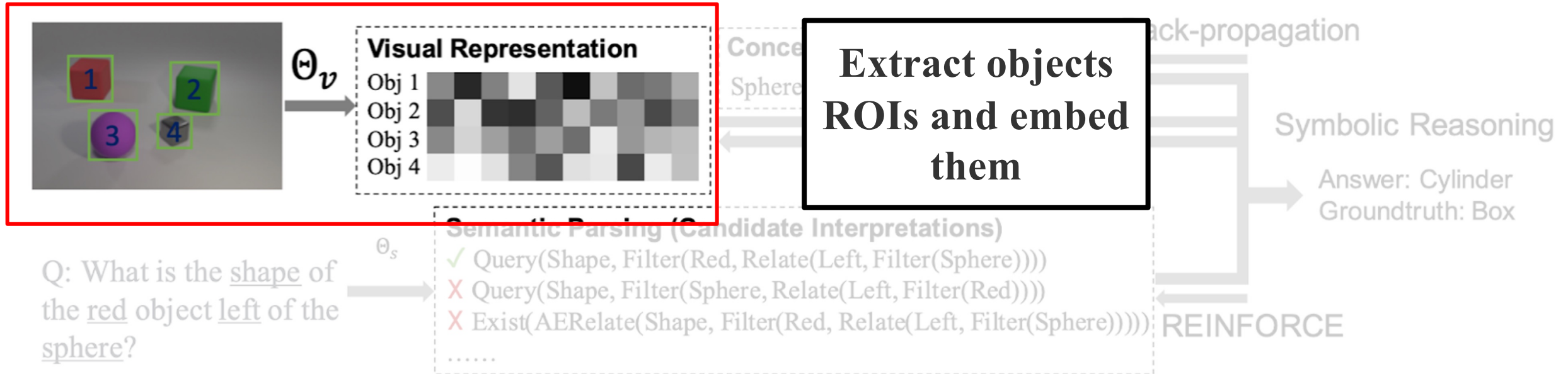


Modular Systems



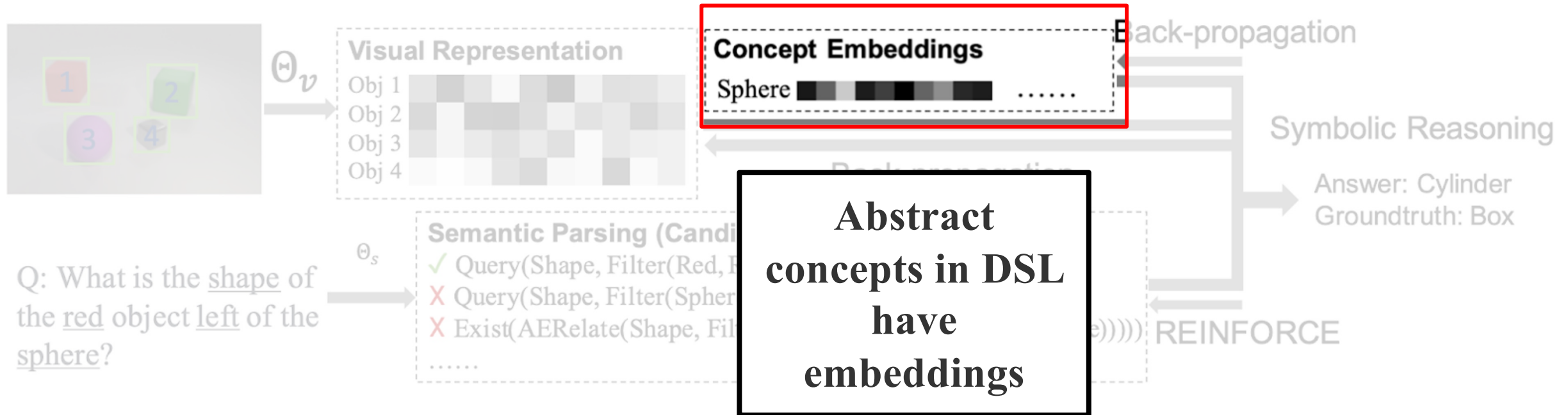


Modular Systems





Modular Systems





Modular Systems

Signature	Implementation
$\text{Scene}() \rightarrow \text{out: ObjectSet}$	$\text{out}_i := 1$
$\text{Filter}(\text{in: ObjectSet}, \text{oc: ObjConcept}) \rightarrow \text{out: ObjectSet}$	$\text{out}_i := \min(\text{in}_i, \text{ObjClassify}(\text{oc})_i)$
$\text{Relate}(\text{in: Object}, \text{rc: RelConcept}) \rightarrow \text{out: ObjectSet}$	$\text{out}_i := \sum_j (\text{in}_j \cdot \text{RelClassify}(\text{rc})_{j,i})$
$\text{AERelate}(\text{in: Object}, \text{a: Attribute}) \rightarrow \text{out: ObjectSet}$	$\text{out}_i := \sum_j (\text{in}_j \cdot \text{AEClassify}(\text{a})_{j,i})$
$\text{Intersection}(\text{in}^{(1)}: \text{ObjectSet}, \text{in}^{(2)}: \text{ObjectSet}) \rightarrow \text{out: ObjectSet}$	$\text{out}_i := \min(\text{in}_i^{(1)}, \text{in}_i^{(2)})$
$\text{Union}(\text{in}^{(1)}: \text{ObjectSet}, \text{in}^{(2)}: \text{ObjectSet}) \rightarrow \text{out: ObjectSet}$	$\text{out}_i := \max(\text{in}_i^{(1)}, \text{in}_i^{(2)})$
$\text{Query}(\text{in: Object}, \text{a: Attribute}) \rightarrow \text{out: ObjConcept}$	$\text{Pr}[\text{out} = \text{oc}] := \sum_i \text{in}_i \cdot \frac{\text{ObjClassify}(\text{oc})_i \cdot b_a^{\text{oc}}}{\sum_{\text{oc}'} \text{ObjClassify}(\text{oc}')_i \cdot b_a^{\text{oc}'}}$
$\text{AEQuery}(\text{in}^{(1)}: \text{Object}, \text{in}^{(2)}: \text{Object}, \text{a: Attribute}) \rightarrow b: \text{Bool}$	$b := \sum_i \sum_j (\text{in}_i^{(1)} \cdot \text{in}_j^{(2)} \cdot \text{AEClassify}(\text{a})_{j,i})$
$\text{Exist}(\text{in: ObjectSet}) \rightarrow b: \text{Bool}$	$b := \max_i \text{in}_i$
$\text{Count}(\text{in: ObjectSet}) \rightarrow i: \text{Integer}$	$i := \sum_i \text{in}_i$
$\text{CLessThan}(\text{in}^{(1)}: \text{ObjectSet}, \text{in}^{(2)}: \text{ObjectSet}) \rightarrow b: \text{Bool}$	$b := \sigma((\sum_i \text{in}_i^{(2)} - \sum_i \text{in}_i^{(1)} - 1 + \gamma_c) / \tau_c)$
$\text{CGreaterThan}(\text{in}^{(1)}: \text{ObjectSet}, \text{in}^{(2)}: \text{ObjectSet}) \rightarrow b: \text{Bool}$	$b := \sigma((\sum_i \text{in}_i^{(1)} - \sum_i \text{in}_i^{(2)} - 1 + \gamma_c) / \tau_c)$
$\text{CEqual}(\text{in}^{(1)}: \text{ObjectSet}, \text{in}^{(2)}: \text{ObjectSet}) \rightarrow b: \text{Bool}$	$b := \sigma((- \sum_i \text{in}_i^{(1)} - \sum_i \text{in}_i^{(2)} + \gamma_c) / (\gamma_c \cdot \tau_c))$

**All operations
deterministic
and pre-defined!**

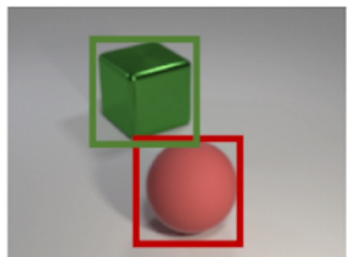
Back-propagation

Symbolic Reasoning

ver: Cylinder
ndtruth: Box



Modular Systems



Object
Detection
→
Feature
Extraction

Obj. 1 
Obj. 2 

a. Perception Module

Parsing
→

Q: Is there any red object?

P: Exist(Filter(red))

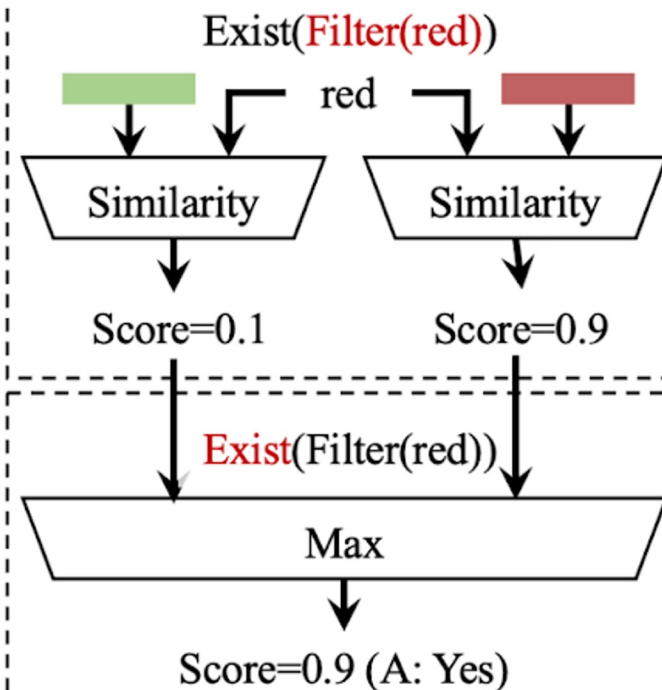
Parsing
→

Q: Do red and green describe the
same property of objects?

P: MetaVerify(red, green, same_kind)

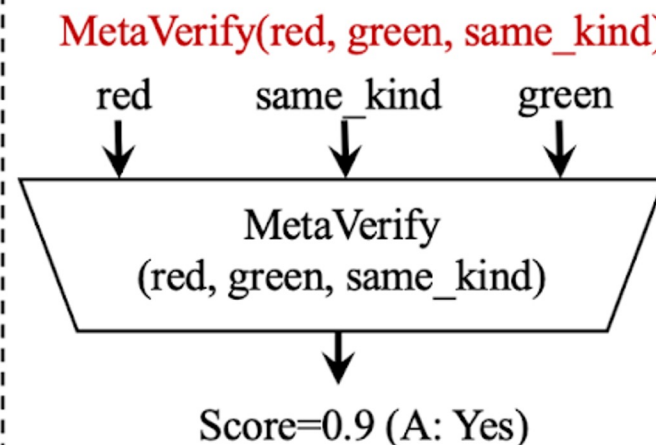
b. Semantic Parsing Module

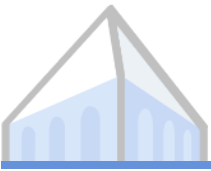
I. Visual Reasoning Question



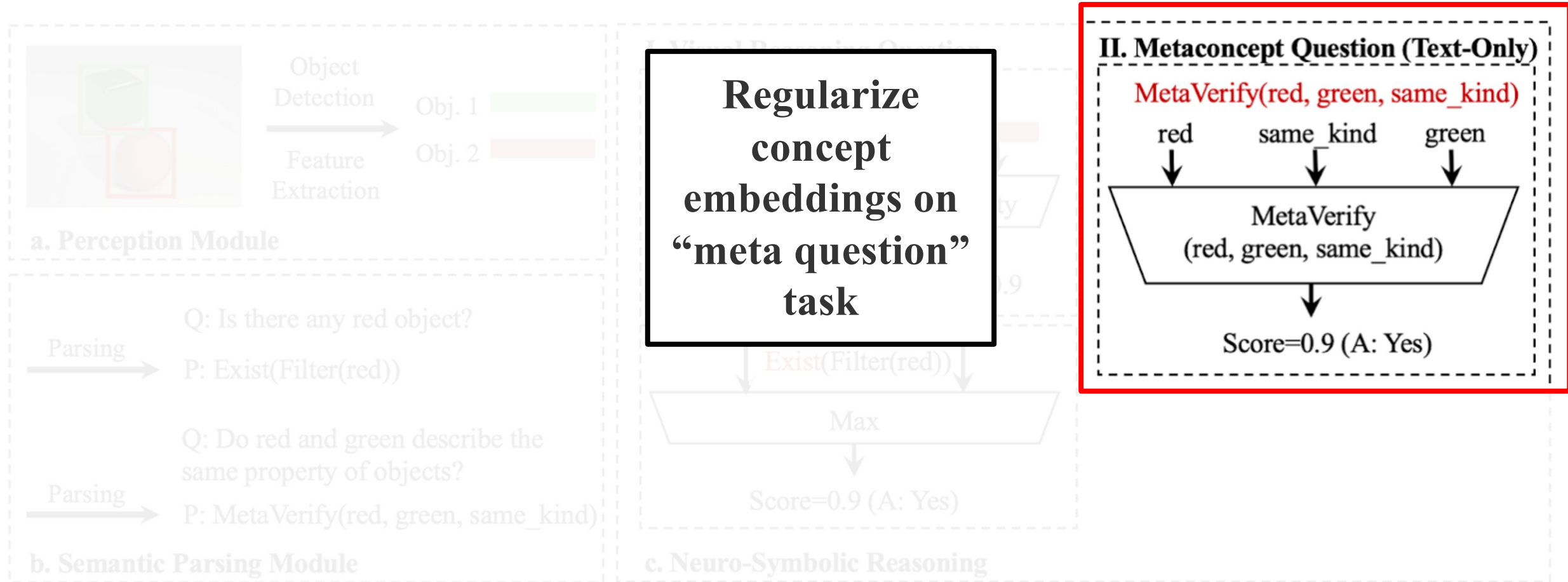
c. Neuro-Symbolic Reasoning

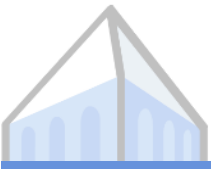
II. Metaconcept Question (Text-Only)





Modular Systems





Modular Systems

“block” == “square”

	GRU-CNN	MAC	NS-CL	VCML
CLEVR	50.0 \pm 0.0	68.7 \pm 3.8	80.2 \pm 3.1	94.1\pm4.6
GQA	50.0 \pm 0.0	49.5 \pm 0.2	49.3 \pm 0.6	50.5\pm0.1

Learning *synonyms* helps zero-shot generalization



Modular Systems



== "purple" + "square"

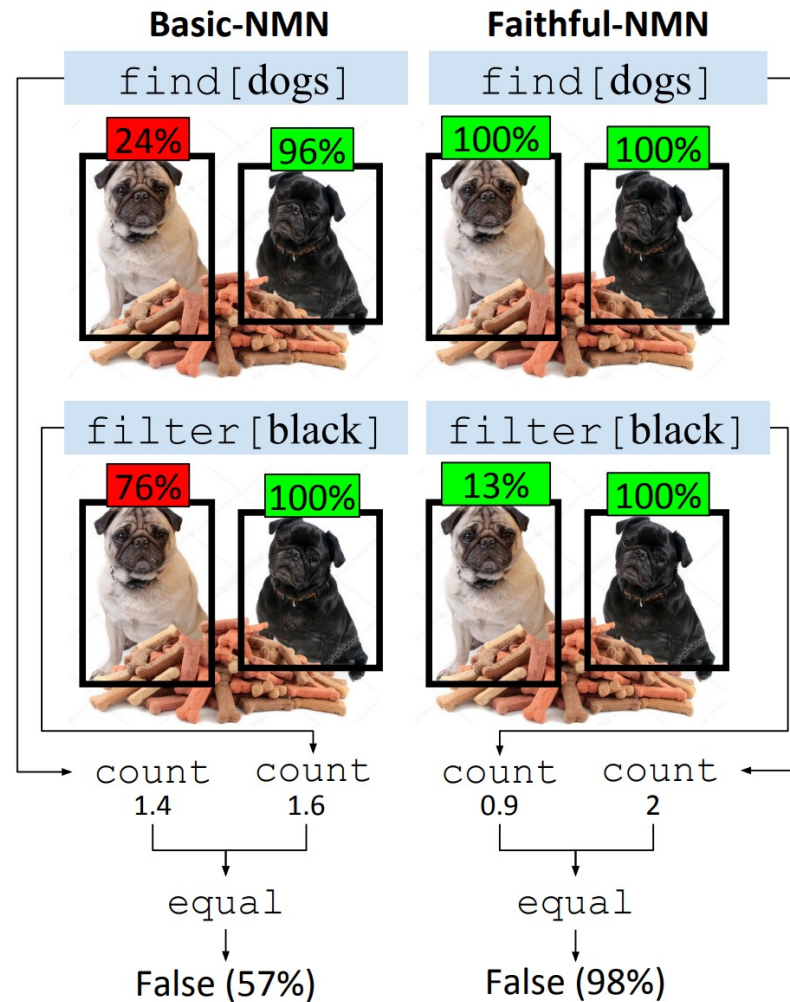
	GRU-CNN	MAC	NS-CL	VCML
CLEVR-200	50.0 \pm 0.0	94.2 \pm 3.3	98.5 \pm 0.3	98.9\pm0.2
CLEVR-20	50.0 \pm 0.0	79.7 \pm 2.6	95.7\pm0.0	95.1 \pm 1.6

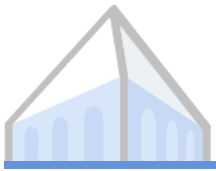
Learning *same kind* helps compositional generalization



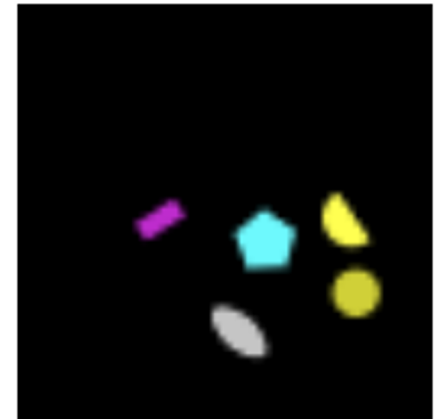
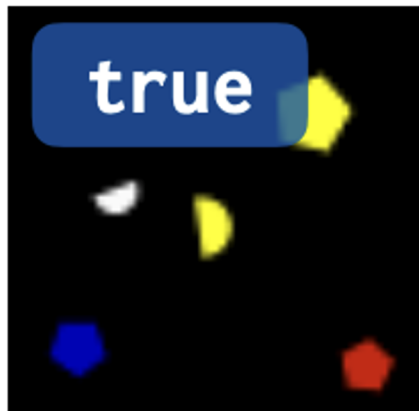
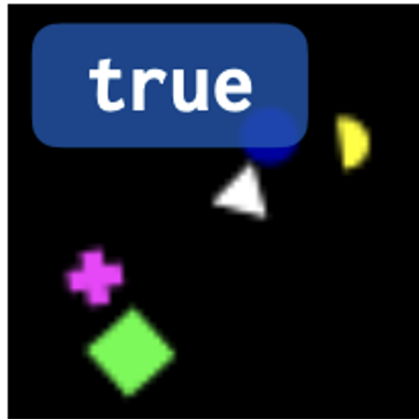
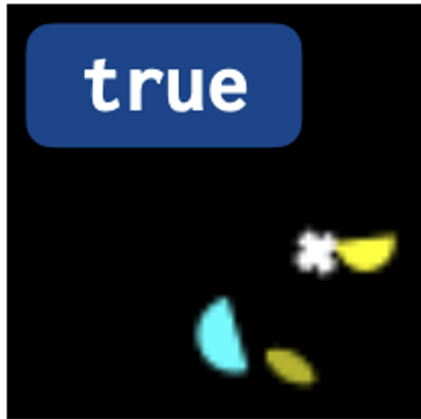
Modular Systems

“All the dogs are black.”



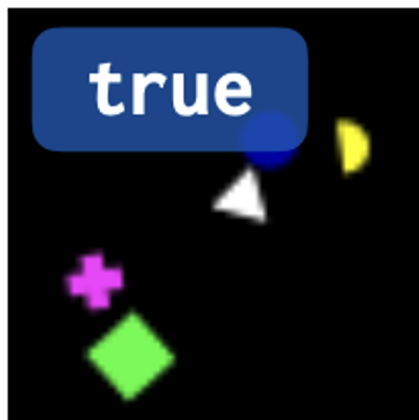
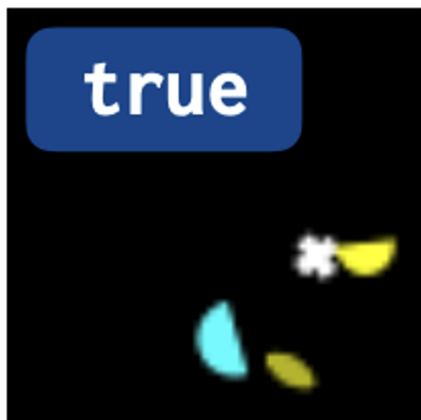


Language as Signal for Abstractions



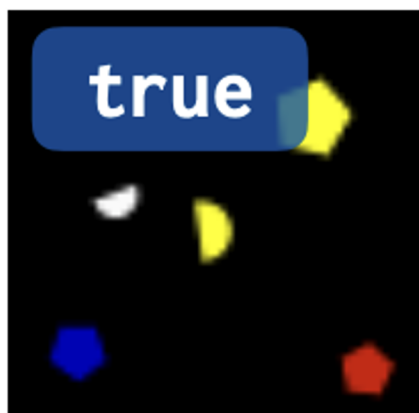


Language as Signal for Abstractions



Available at Training

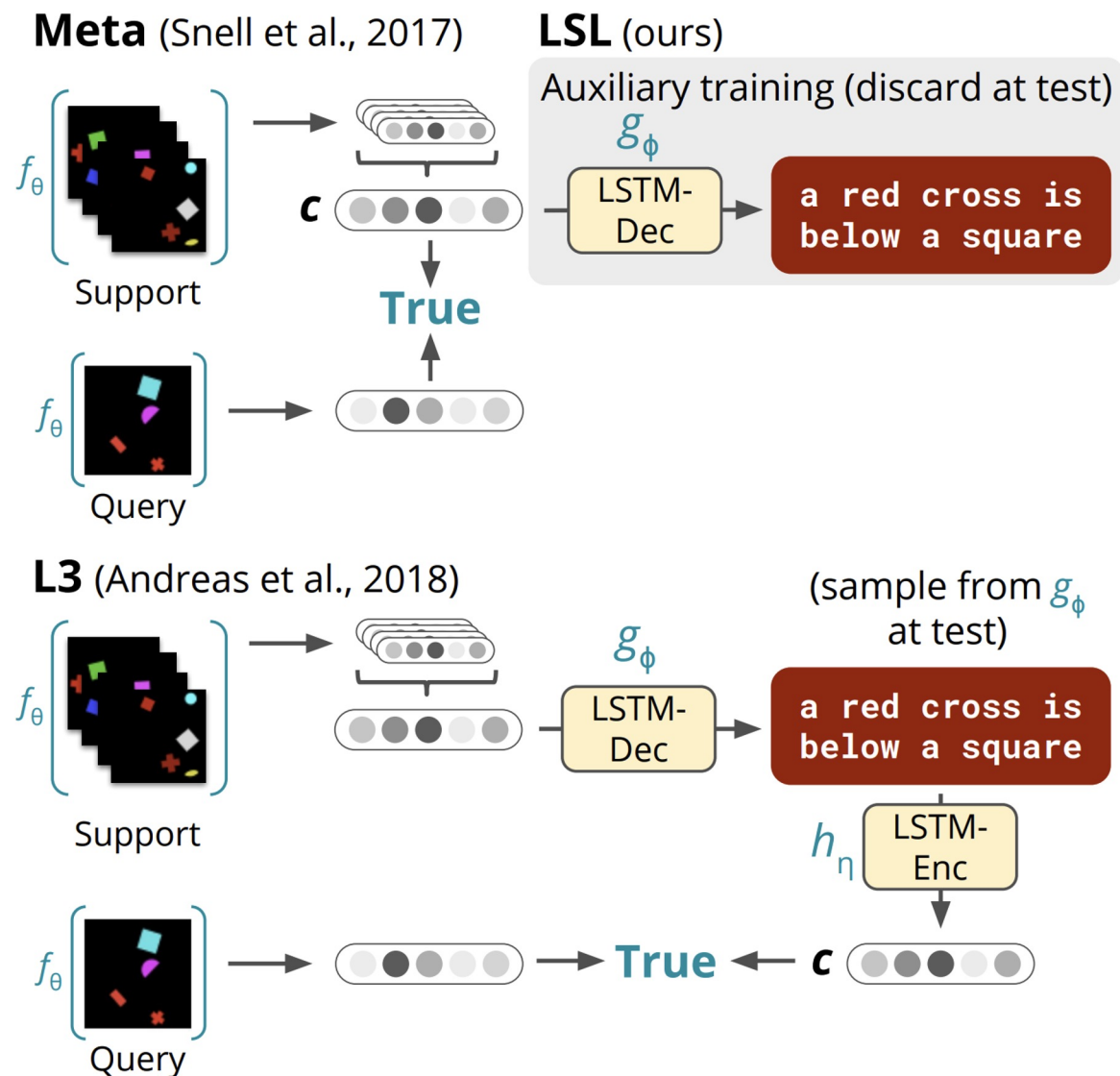
a white shape is
left of a yellow
semicircle



true

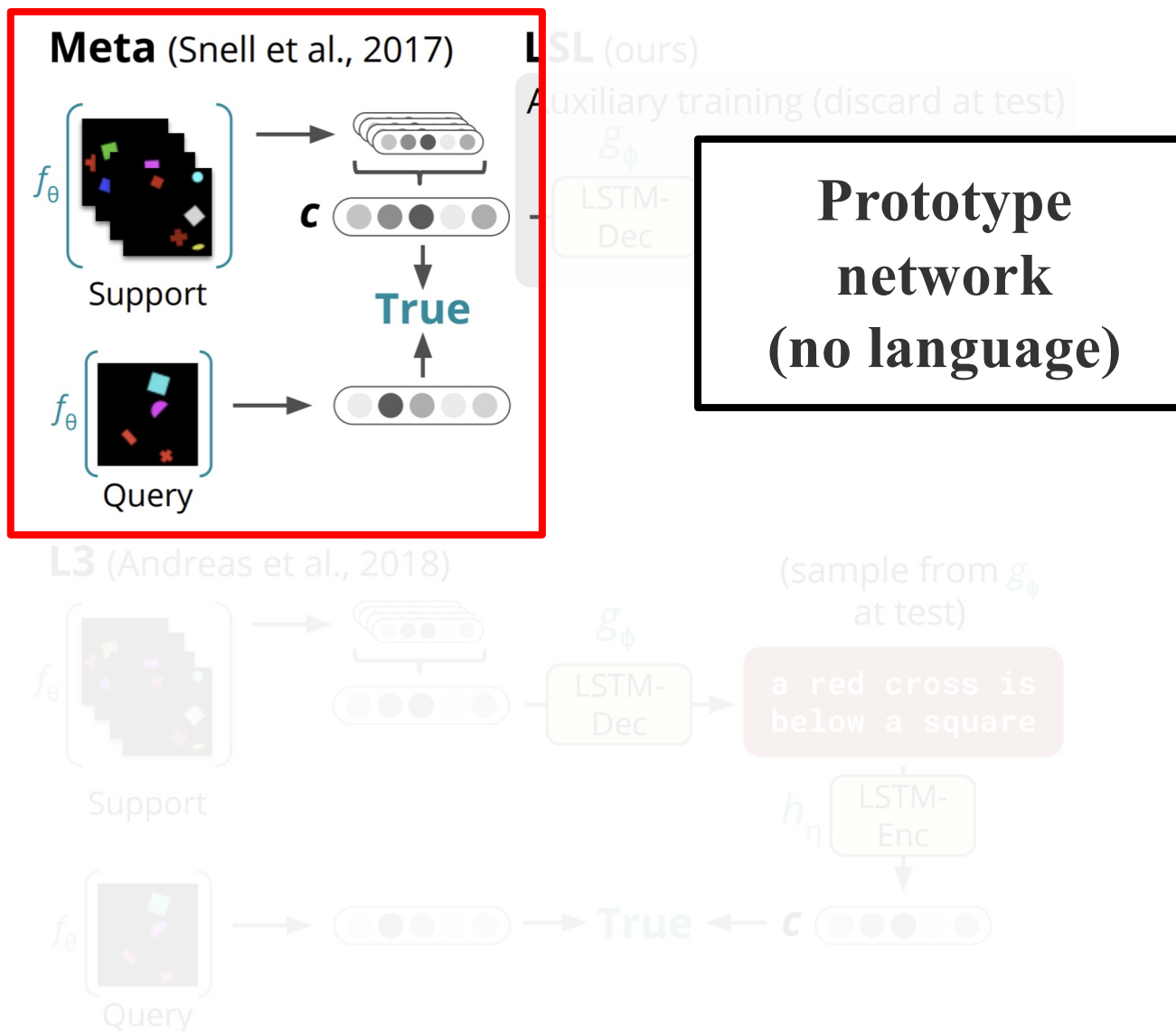


Language as Signal for Abstractions



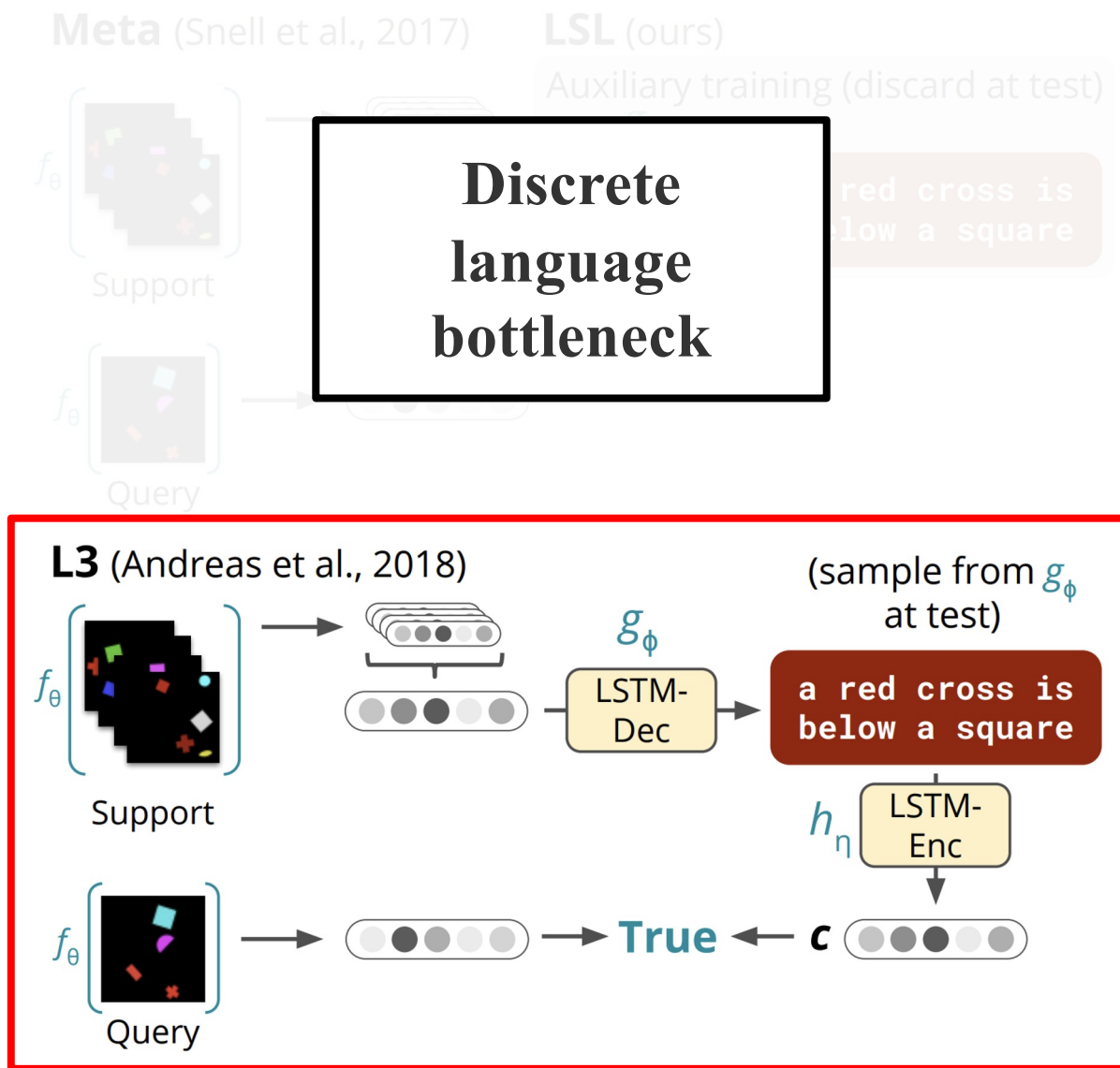


Language as Signal for Abstractions



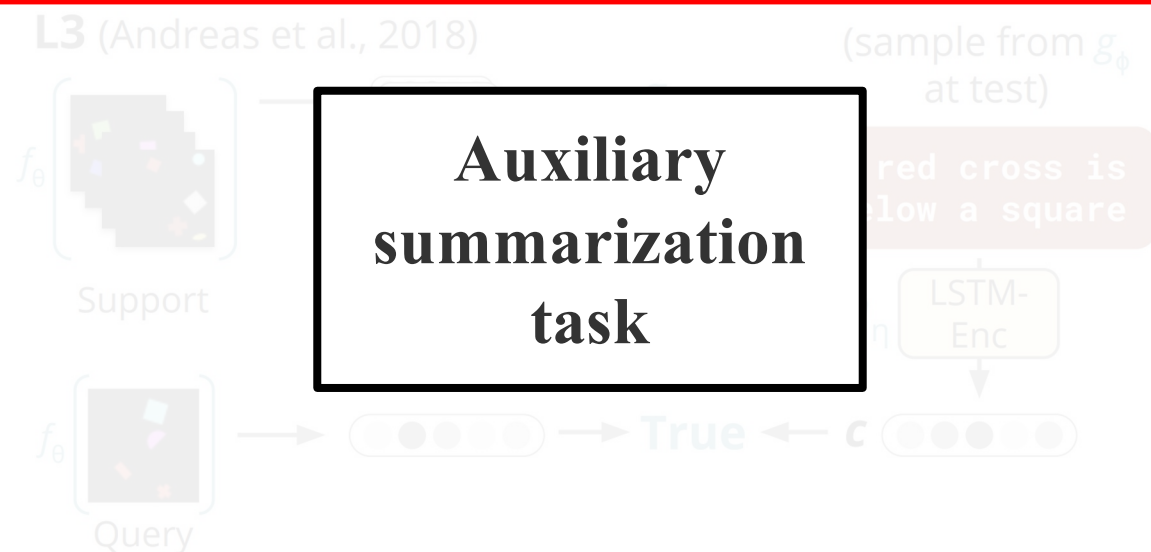
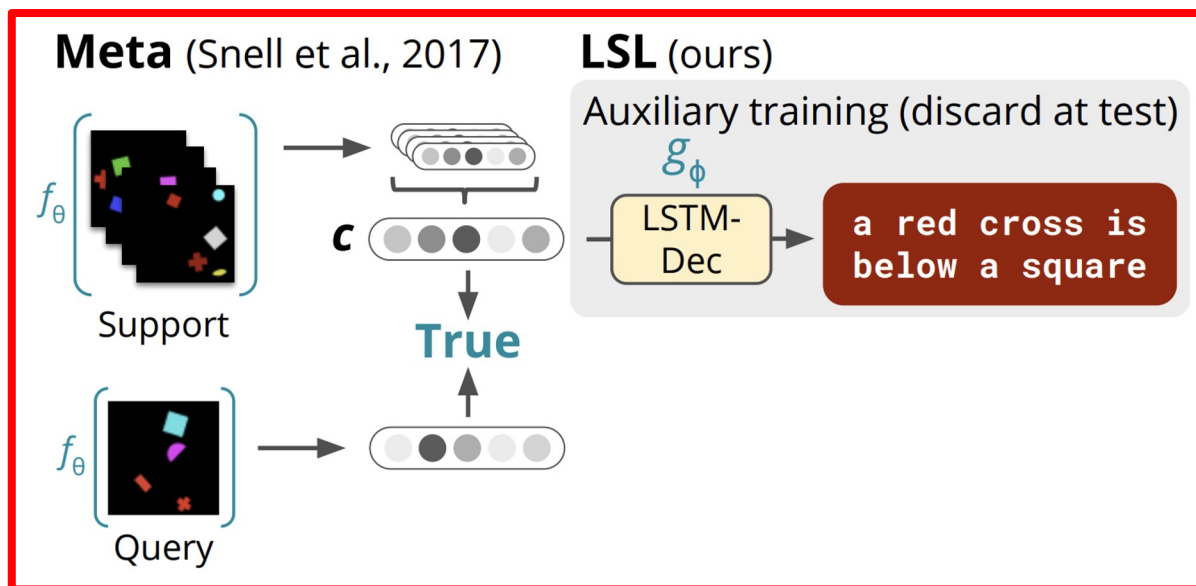


Language as Signal for Abstractions





Language as Signal for Abstractions



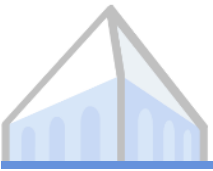


Language as Signal for Abstractions

Test Set Accuracy

	ShapeWorld	Birds ($D = 20$)
--	------------	--------------------

Meta	60.59 \pm 1.07	57.97 \pm 0.96
L3	66.60 \pm 1.18	53.96 \pm 1.06
LSL	67.29 \pm 1.03	61.24 \pm 0.96



Top-Down Takeaways

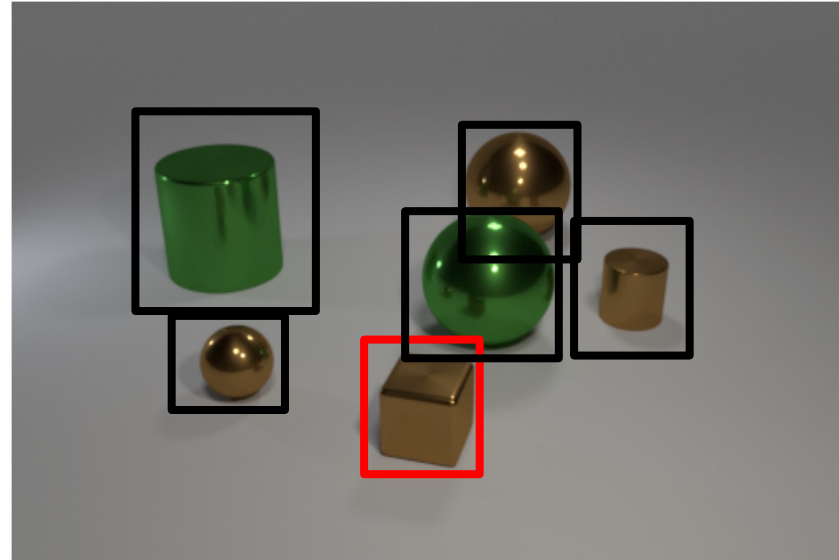
- Language provides labels for supervised learning of perceptual systems.
- Can provide powerful inductive biases in computational structure *if known*.
- Serves as signal for useful perceptual abstractions to learn either as bottleneck or auxiliary signal.

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)



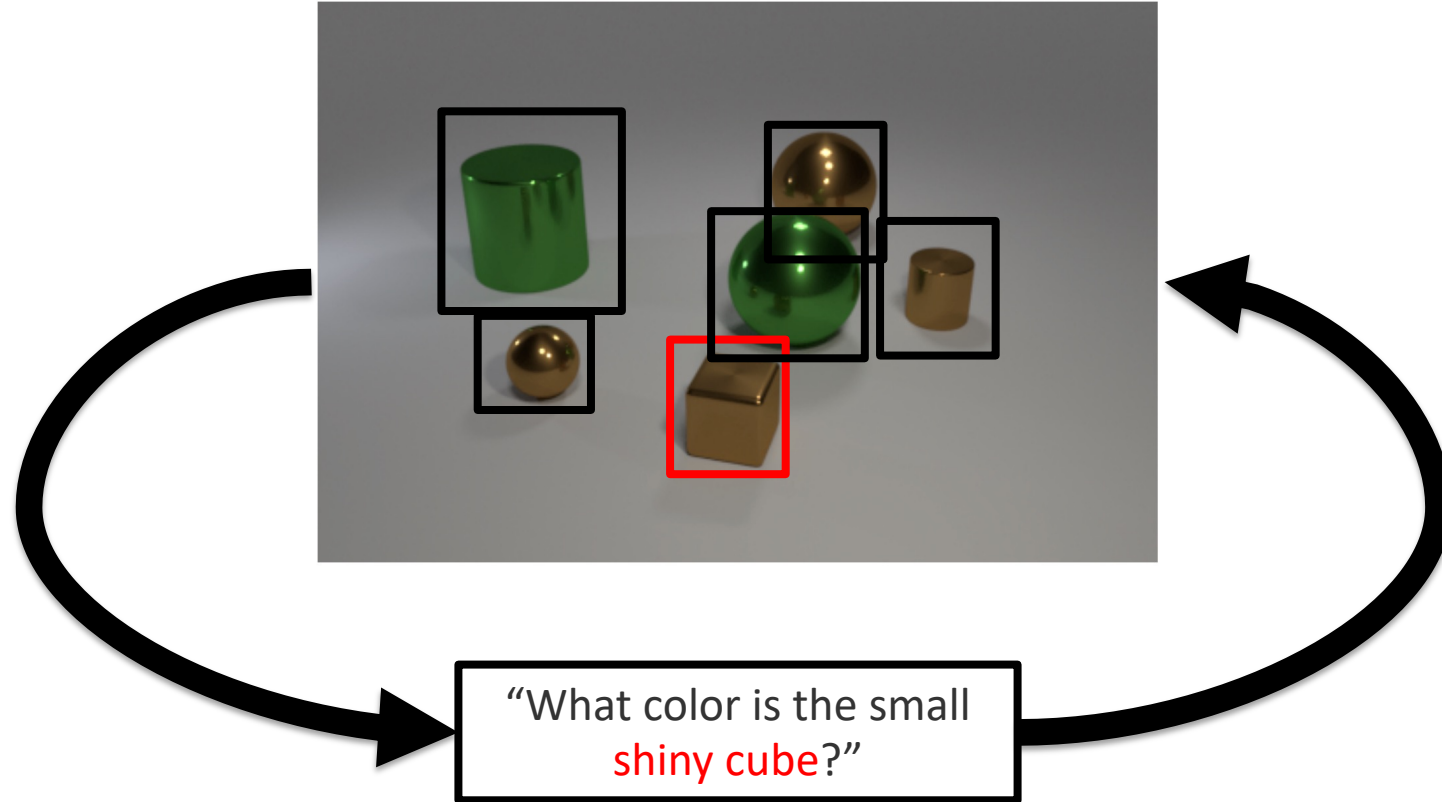
Bottom-Up & Top-Down Reasoning



“What color is the small
shiny cube?”



Bottom-Up & Top-Down Reasoning





Galatea of the Spheres, Salvador Dali 1952

Extra Slides



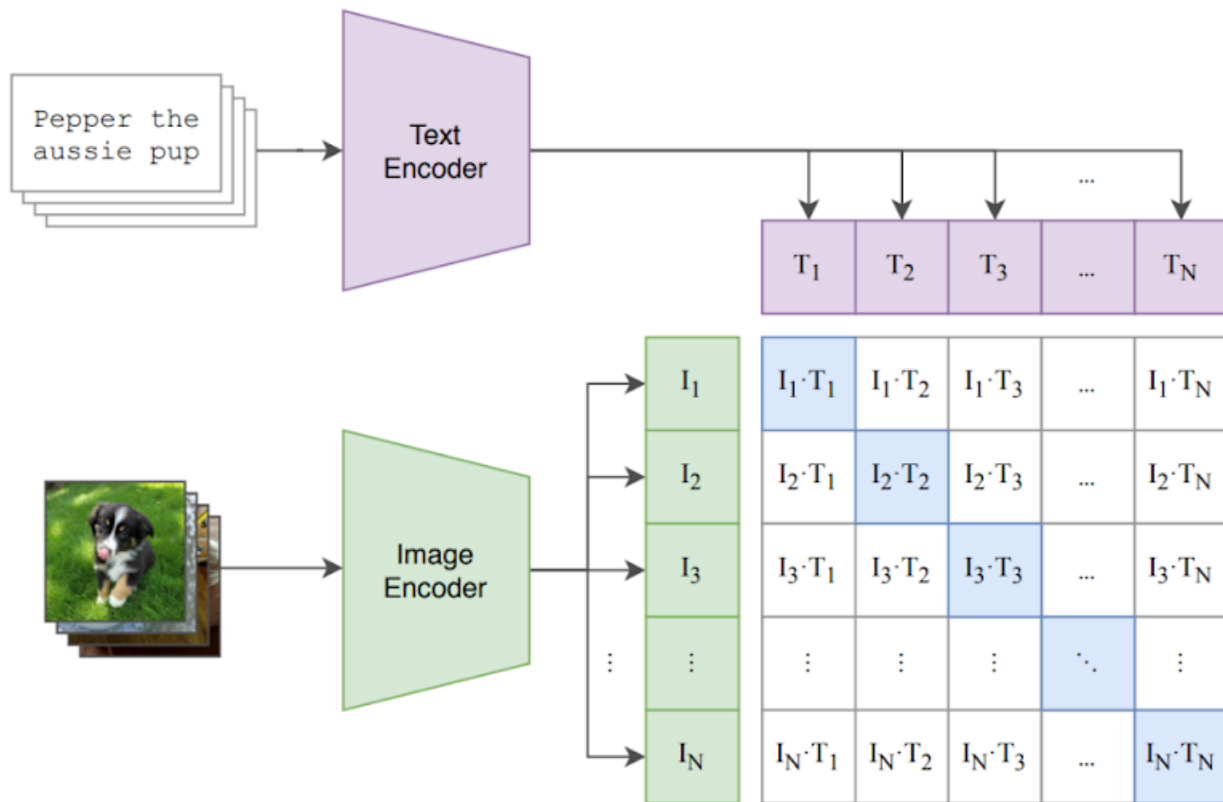
Open-Set Models

Models which leverage the open-vocabulary of language to enjoy a practically open set of labels!

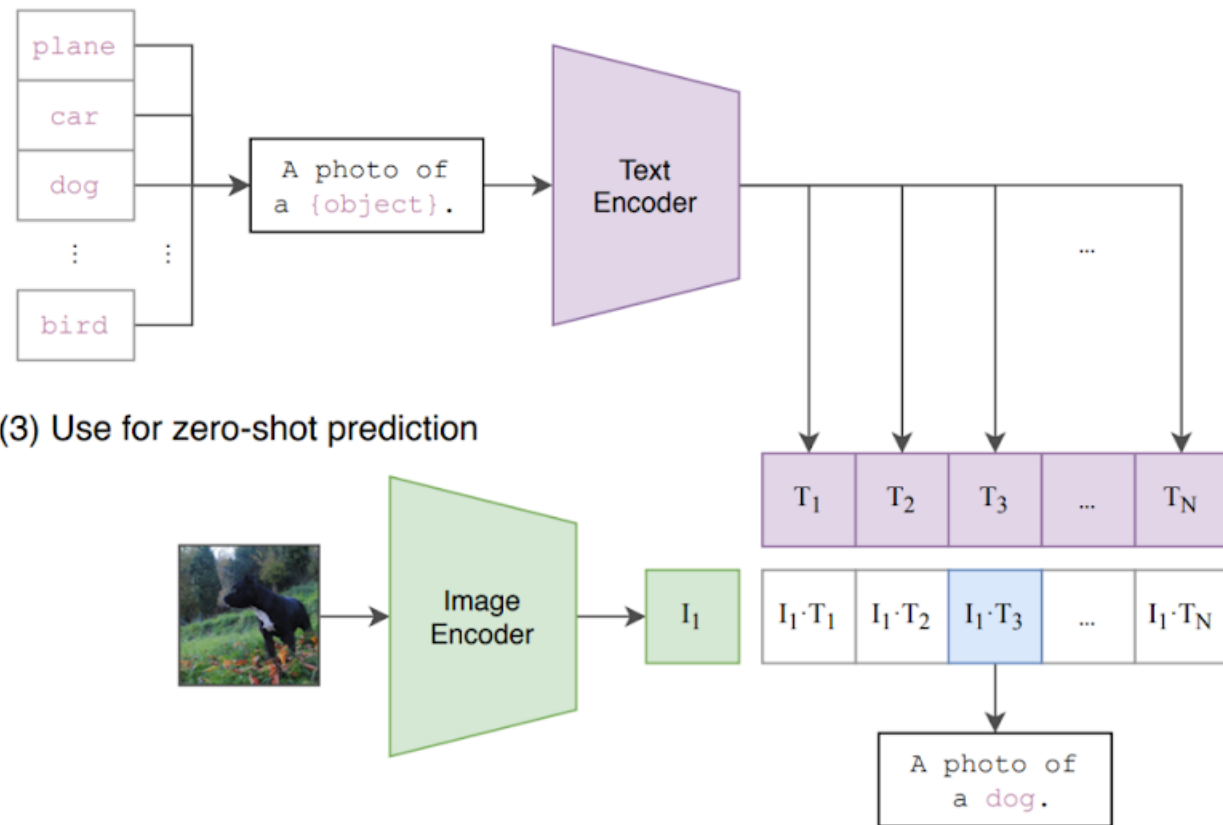


Open-Set Models

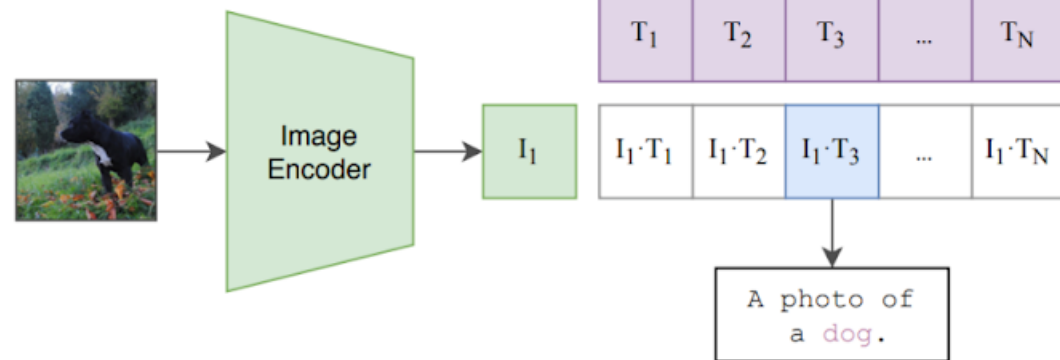
(1) Contrastive pre-training



(2) Create dataset classifier from label text

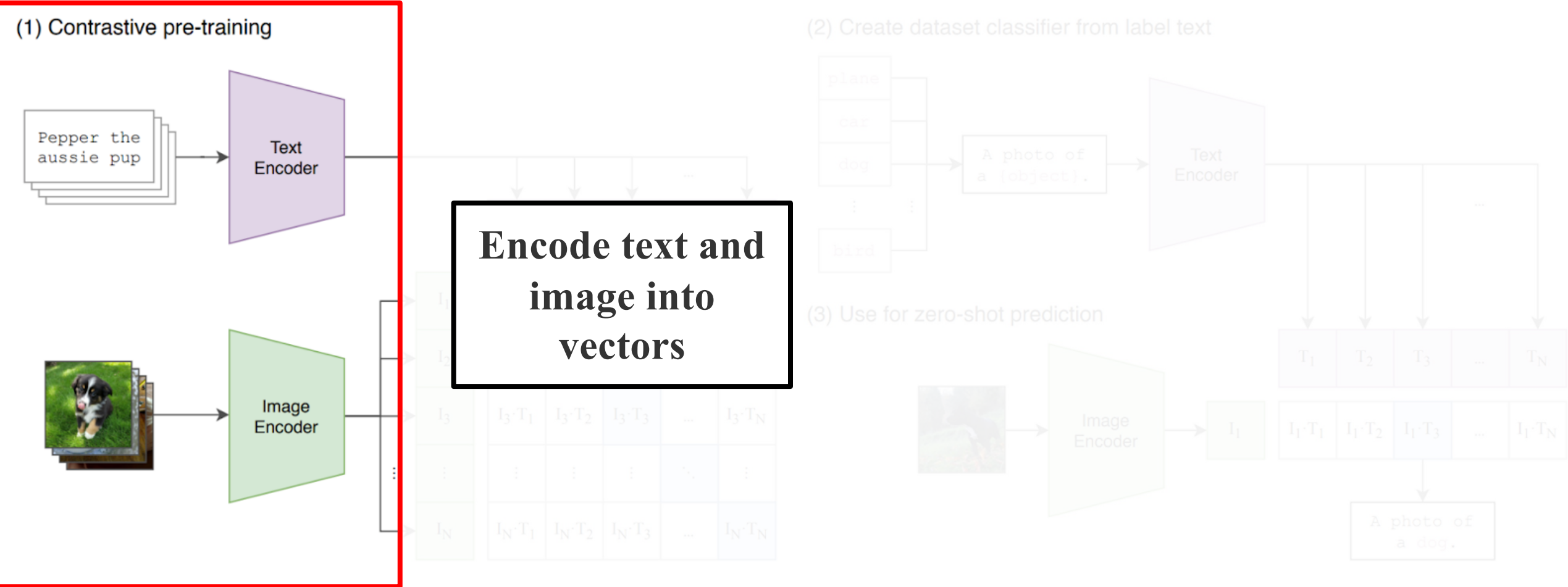


(3) Use for zero-shot prediction





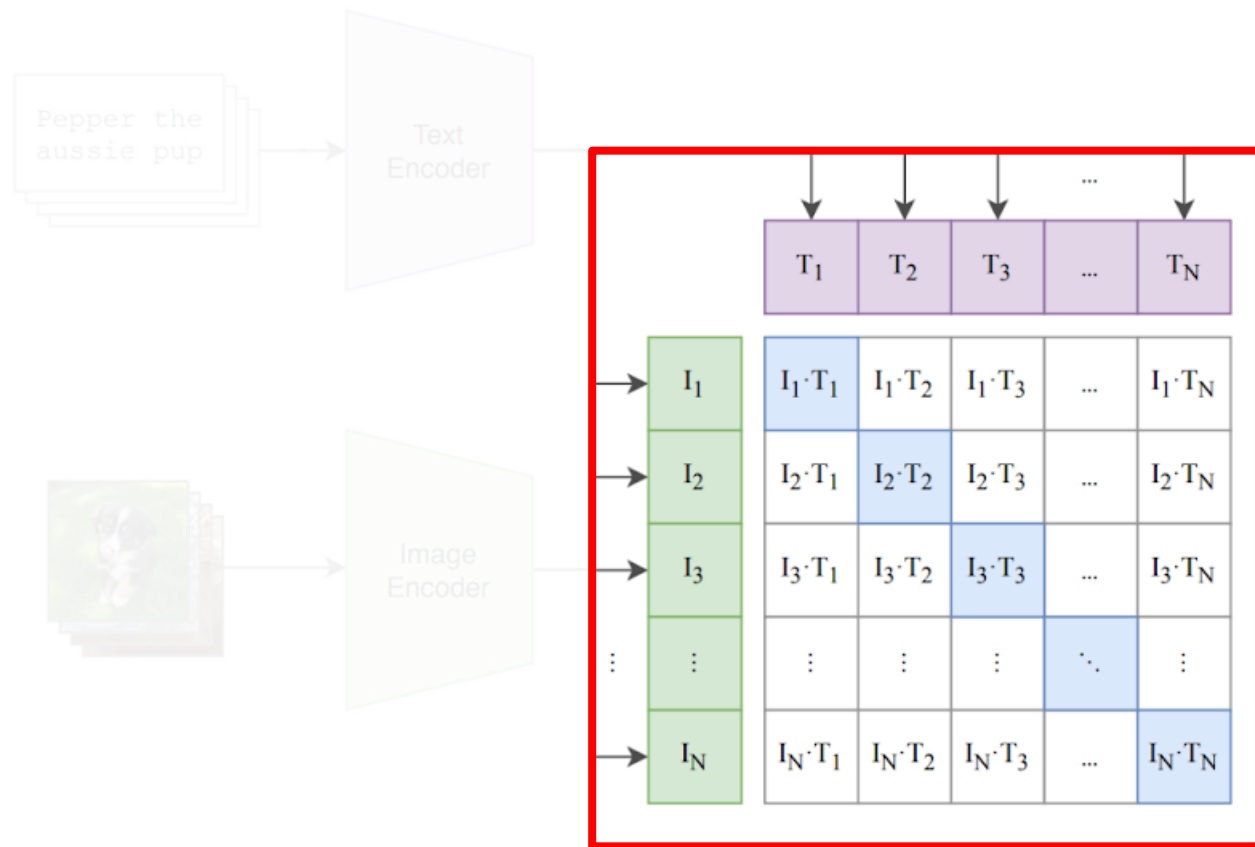
Open-Set Models



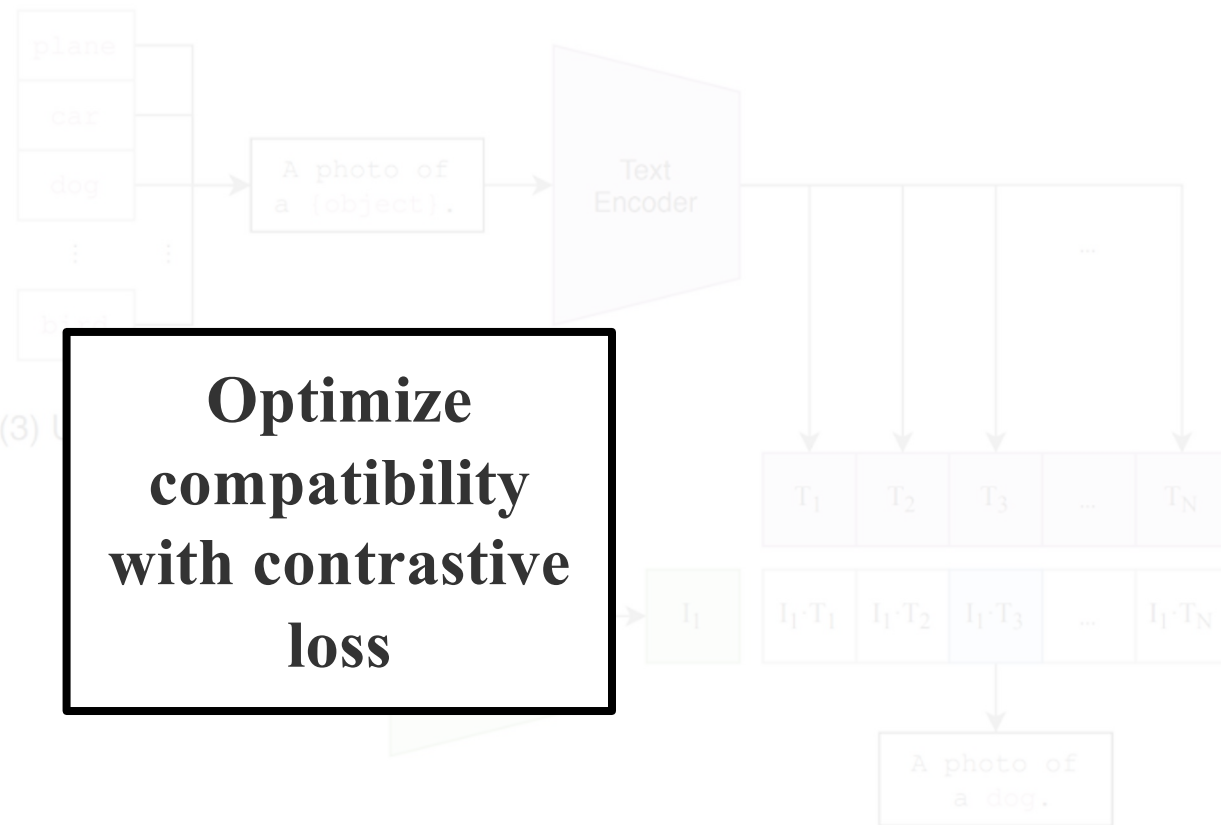


Open-Set Models

(1) Contrastive pre-training



(2) Create dataset classifier from label text

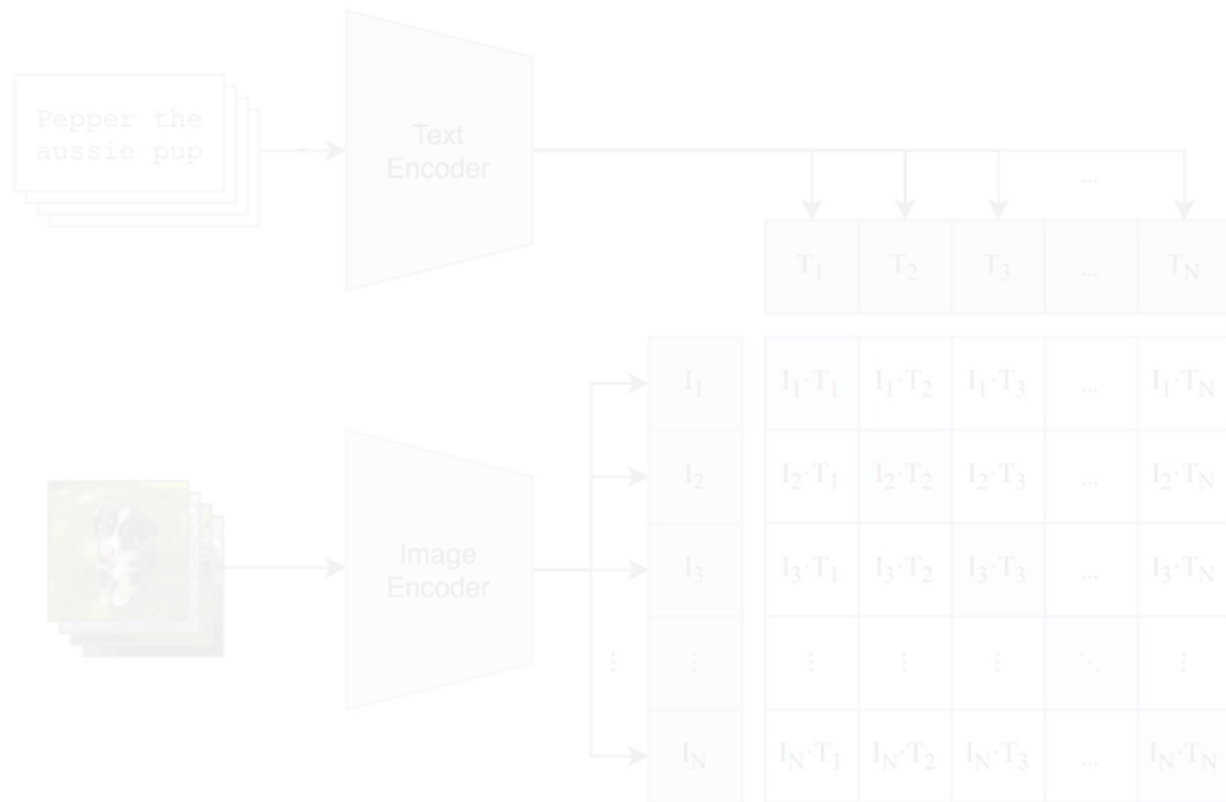


**Optimize
compatibility
with contrastive
loss**

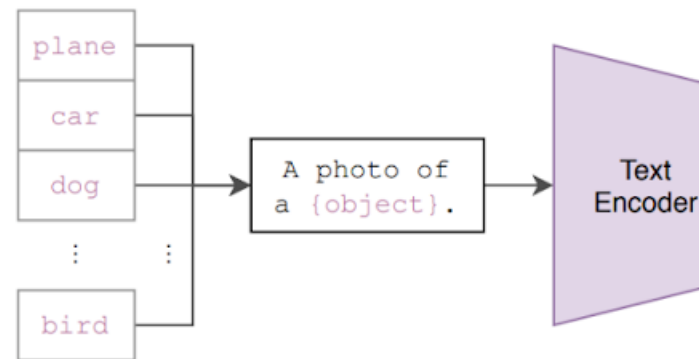


Open-Set Models

(1) Contrastive pre-training

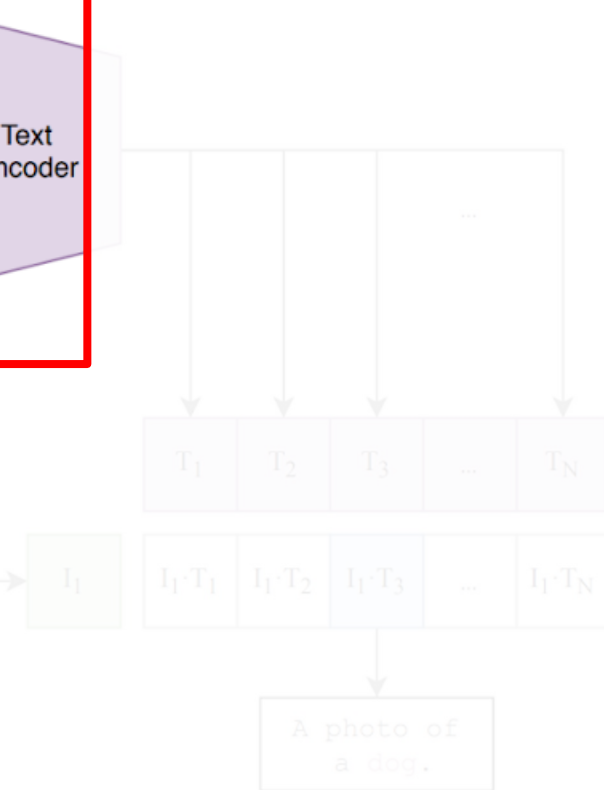


(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

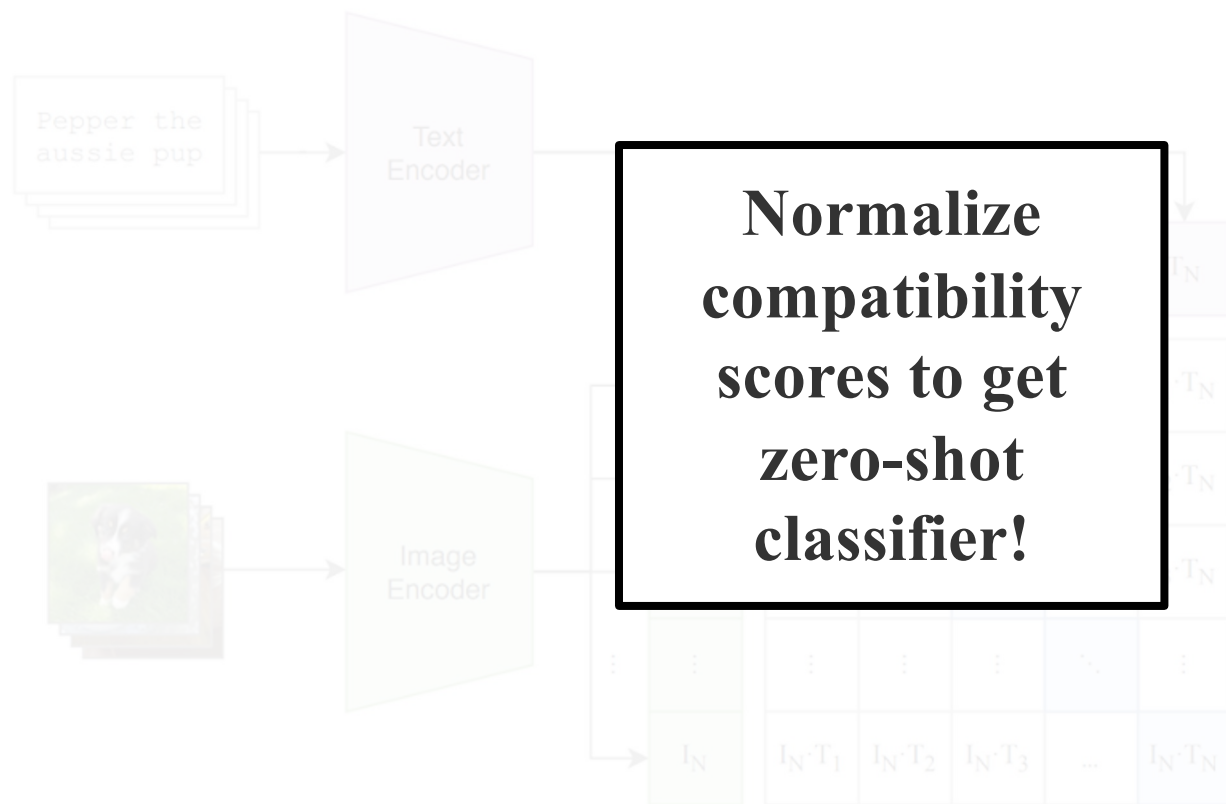
Classification dataset created with templated prompts



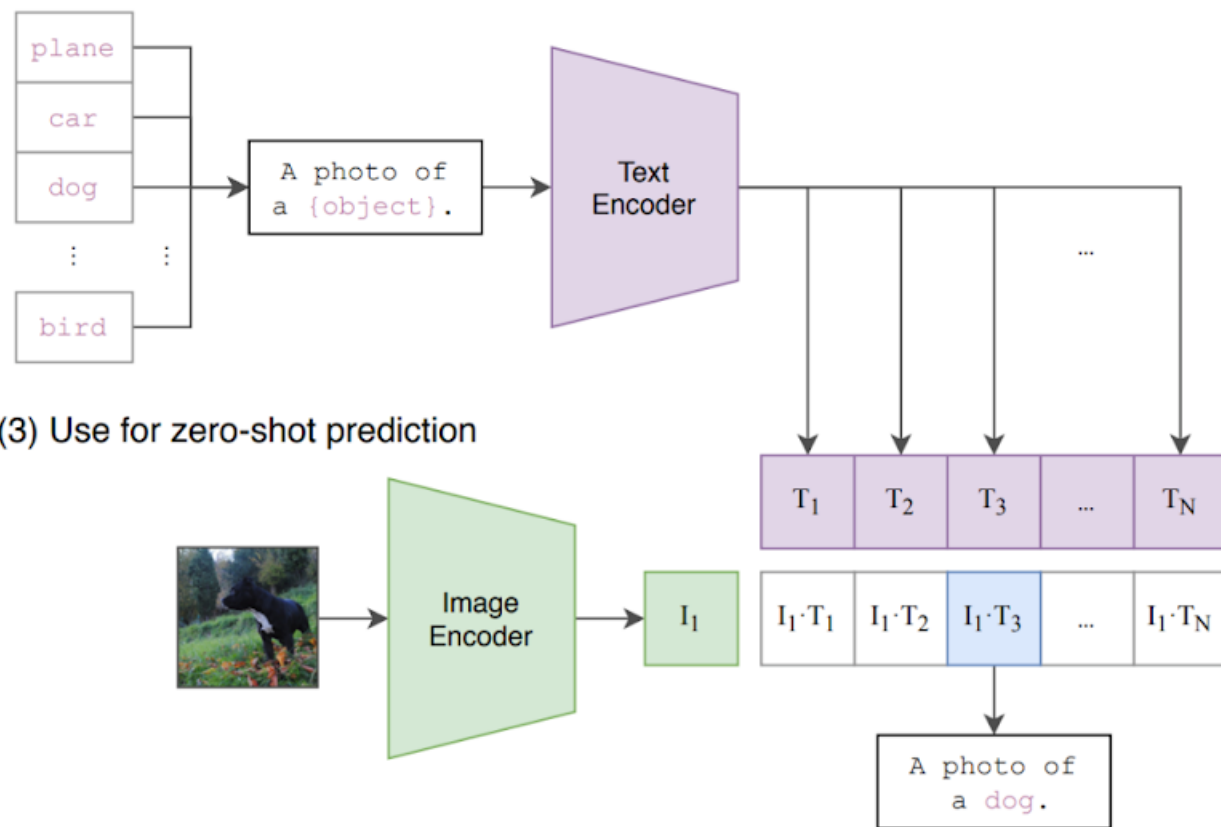


Open-Set Models

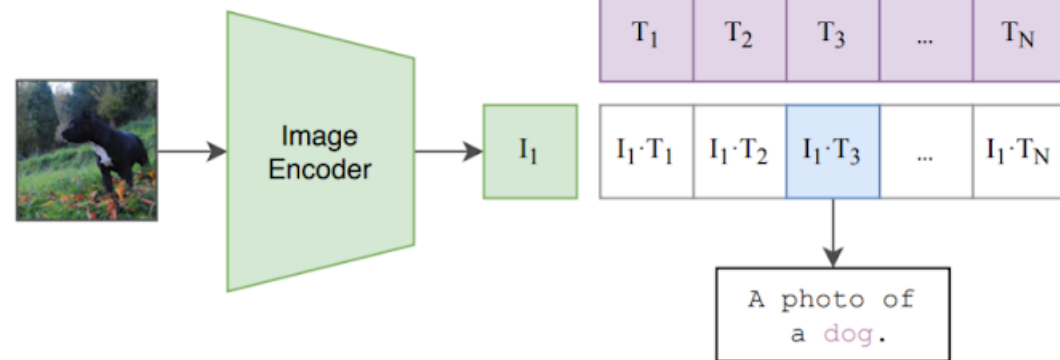
(1) Contrastive pre-training



(2) Create dataset classifier from label text

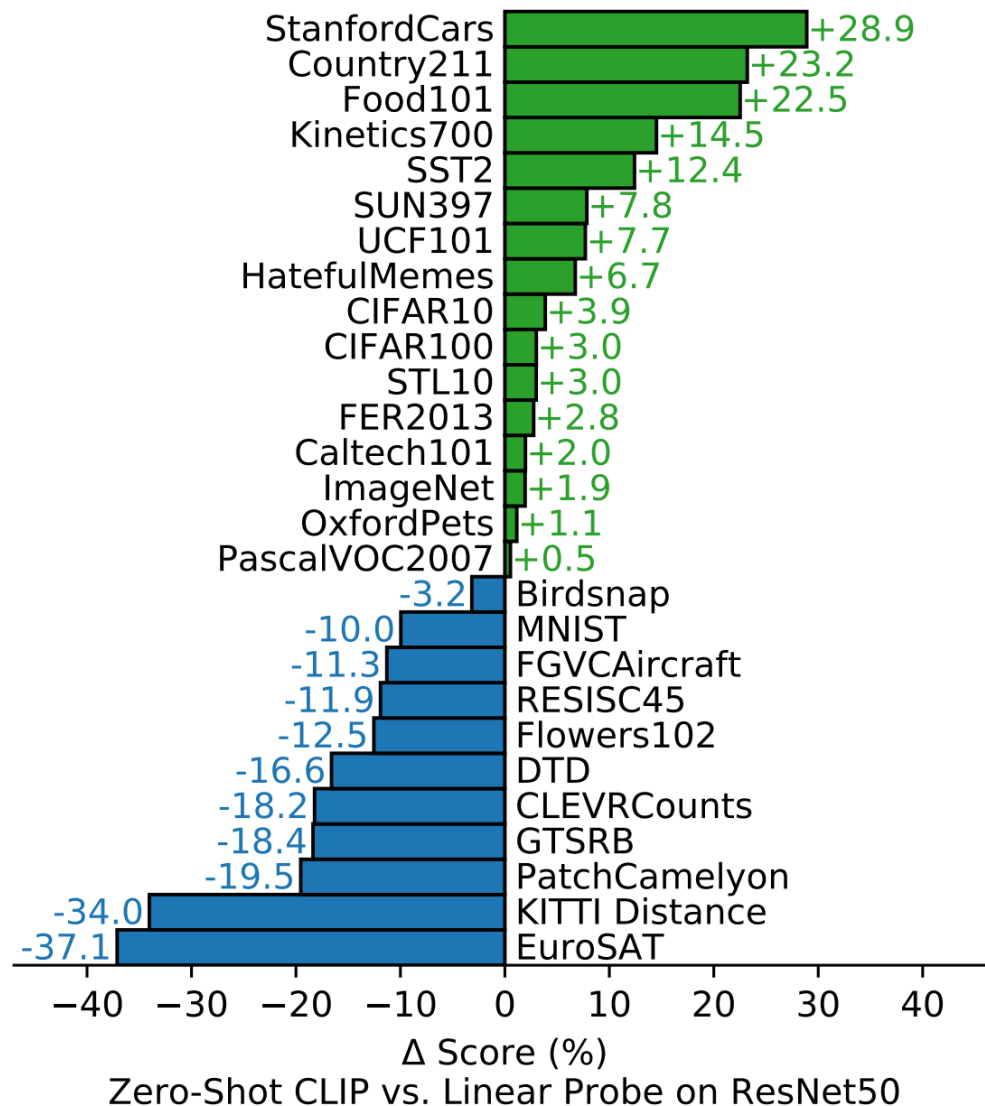


(3) Use for zero-shot prediction



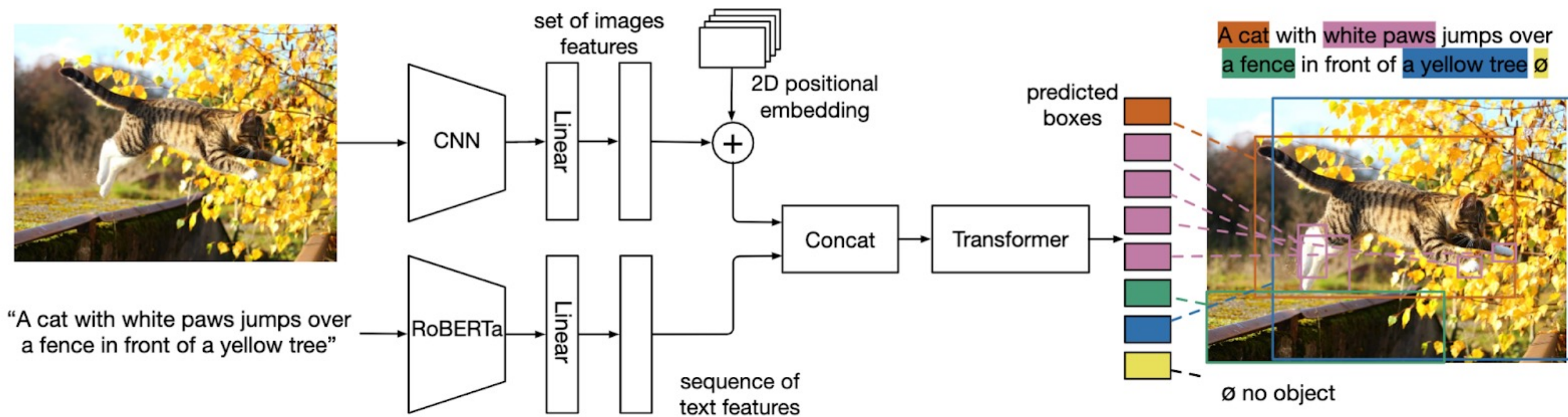


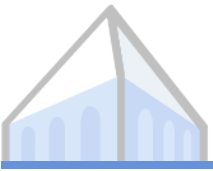
Open-Set Models



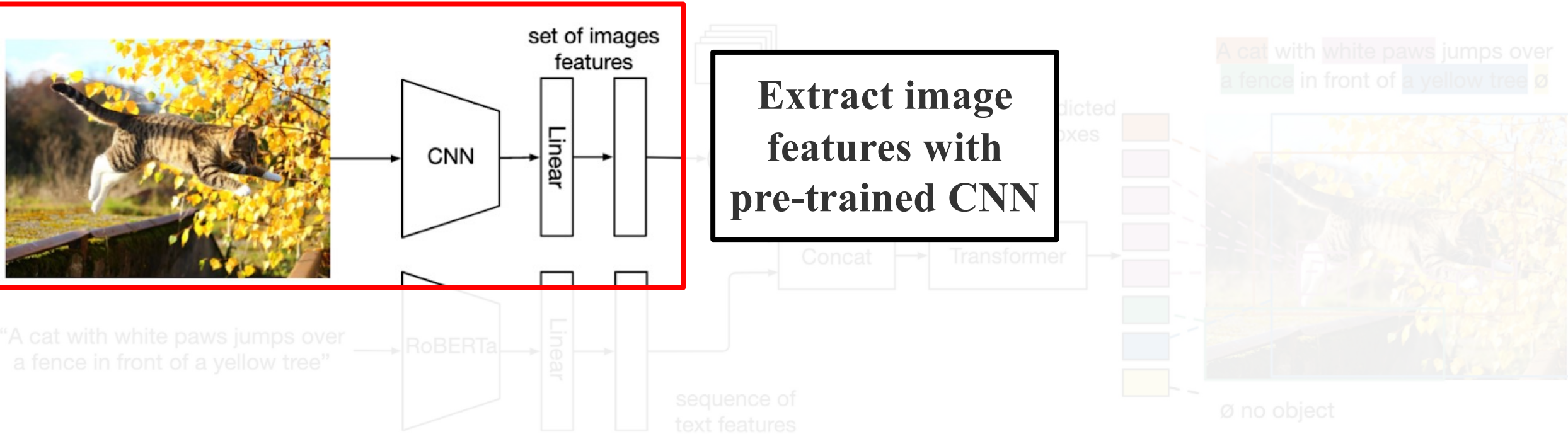


Open-Set Models



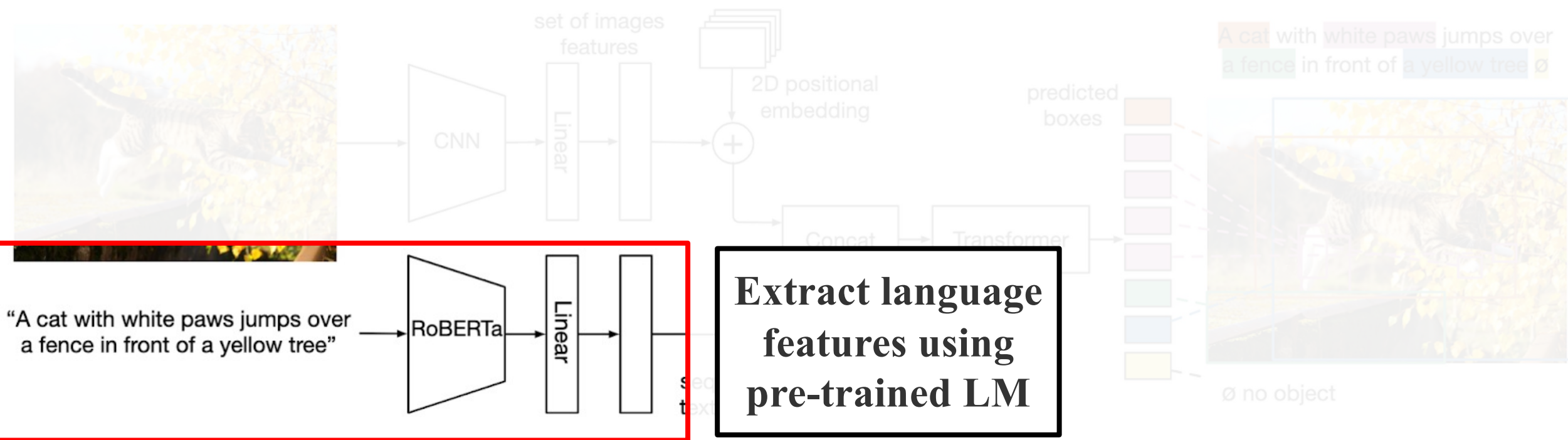


Open-Set Models



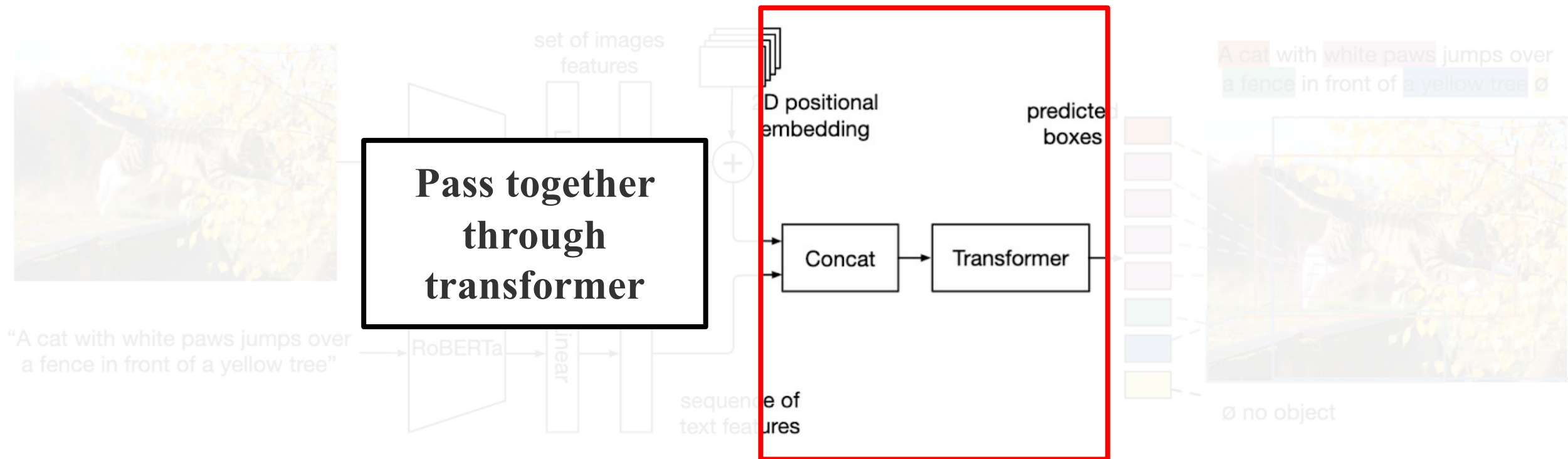


Open-Set Models



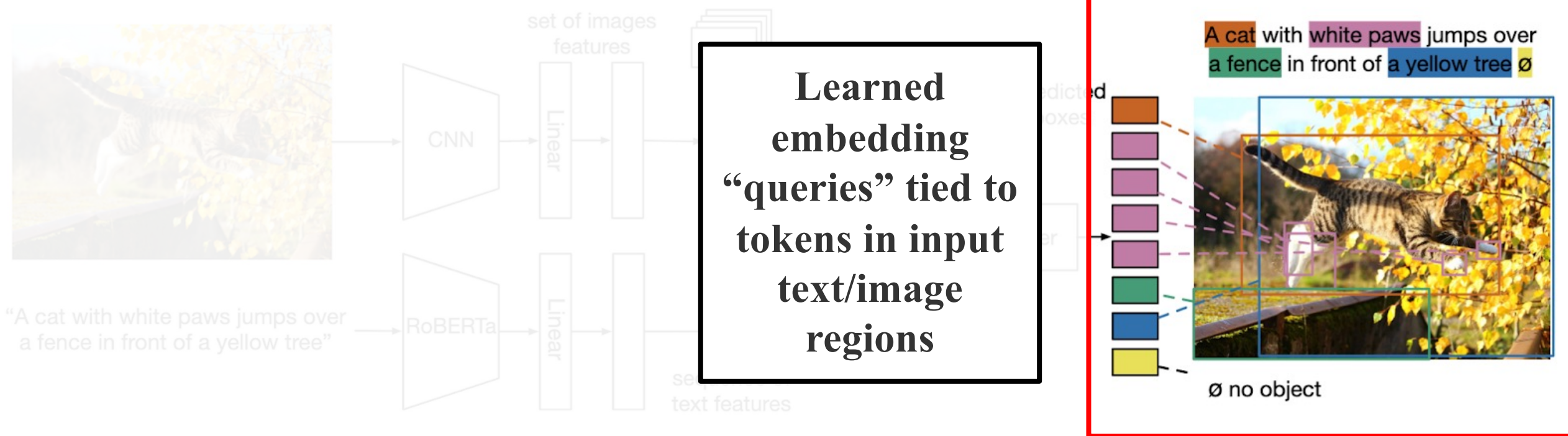


Open-Set Models





Open-Set Models





Open-Set Models

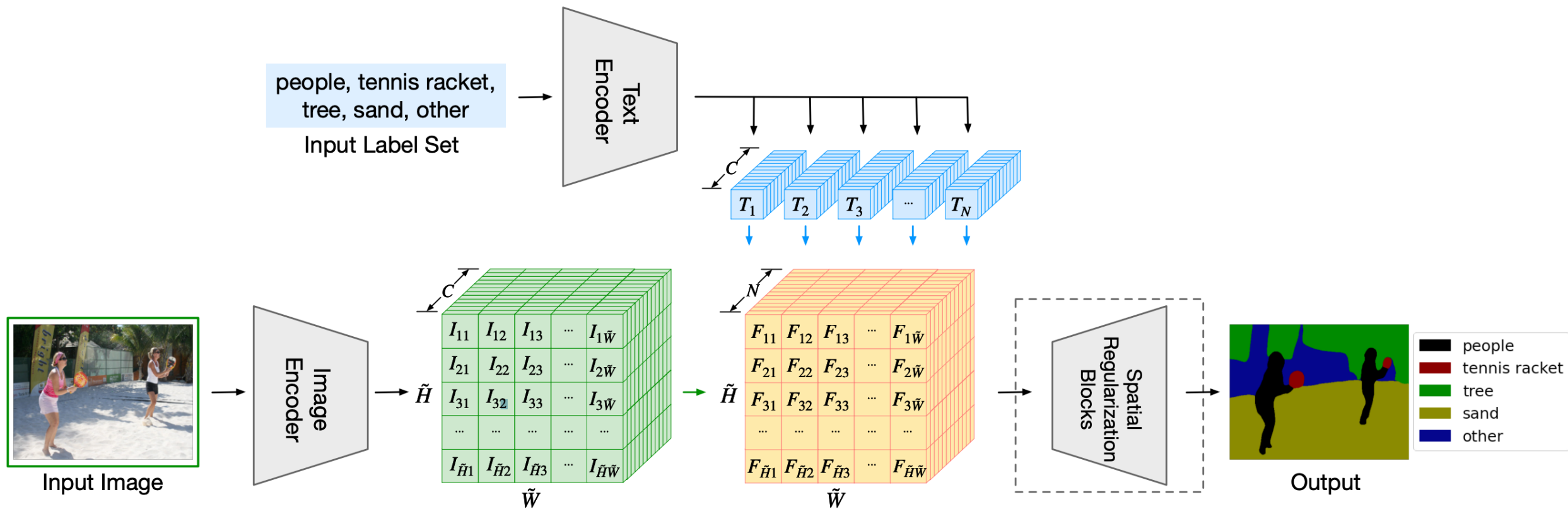


(a) “one small boy climbing a pole with the help of another boy on the ground”

(b) “A man talking on his cellphone next to a jewelry store”

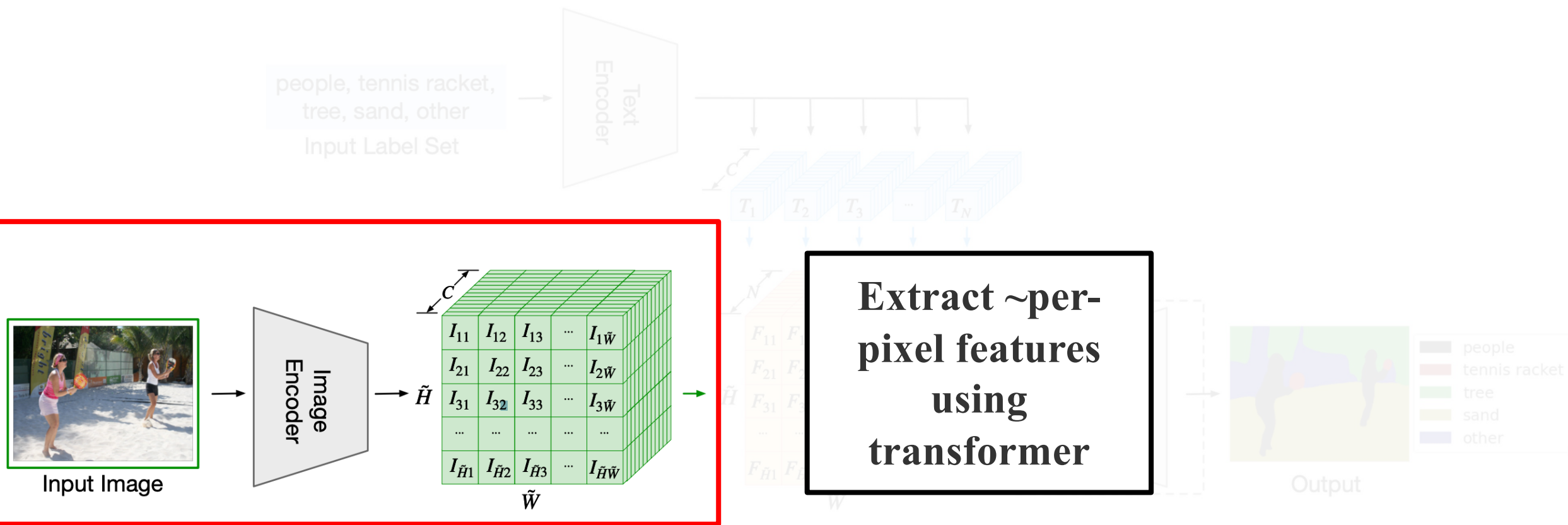


Open-Set Models



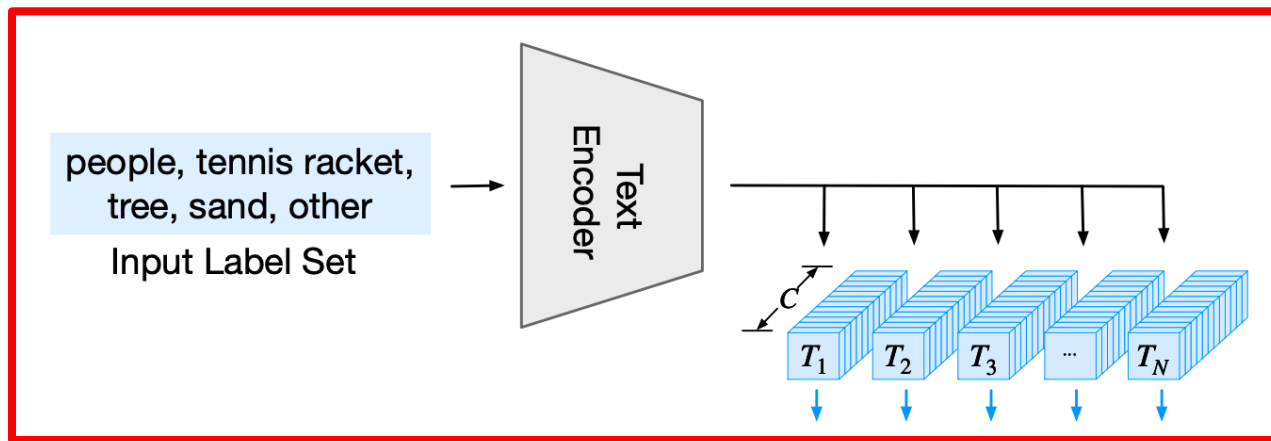


Open-Set Models



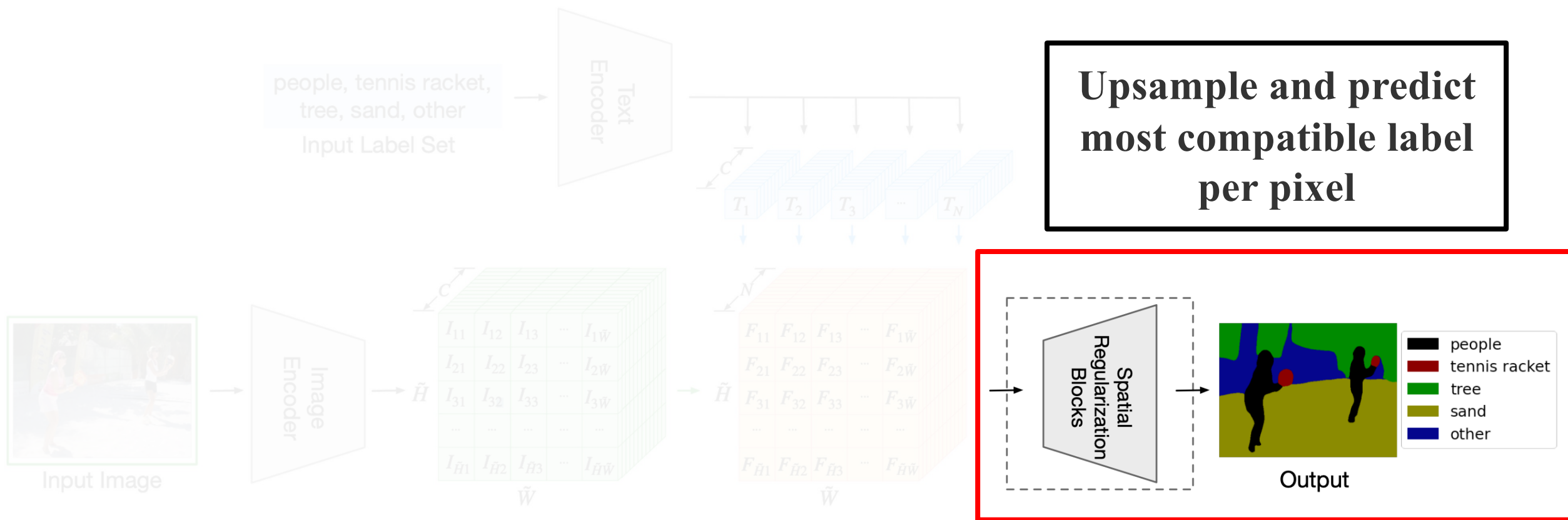


Open-Set Models



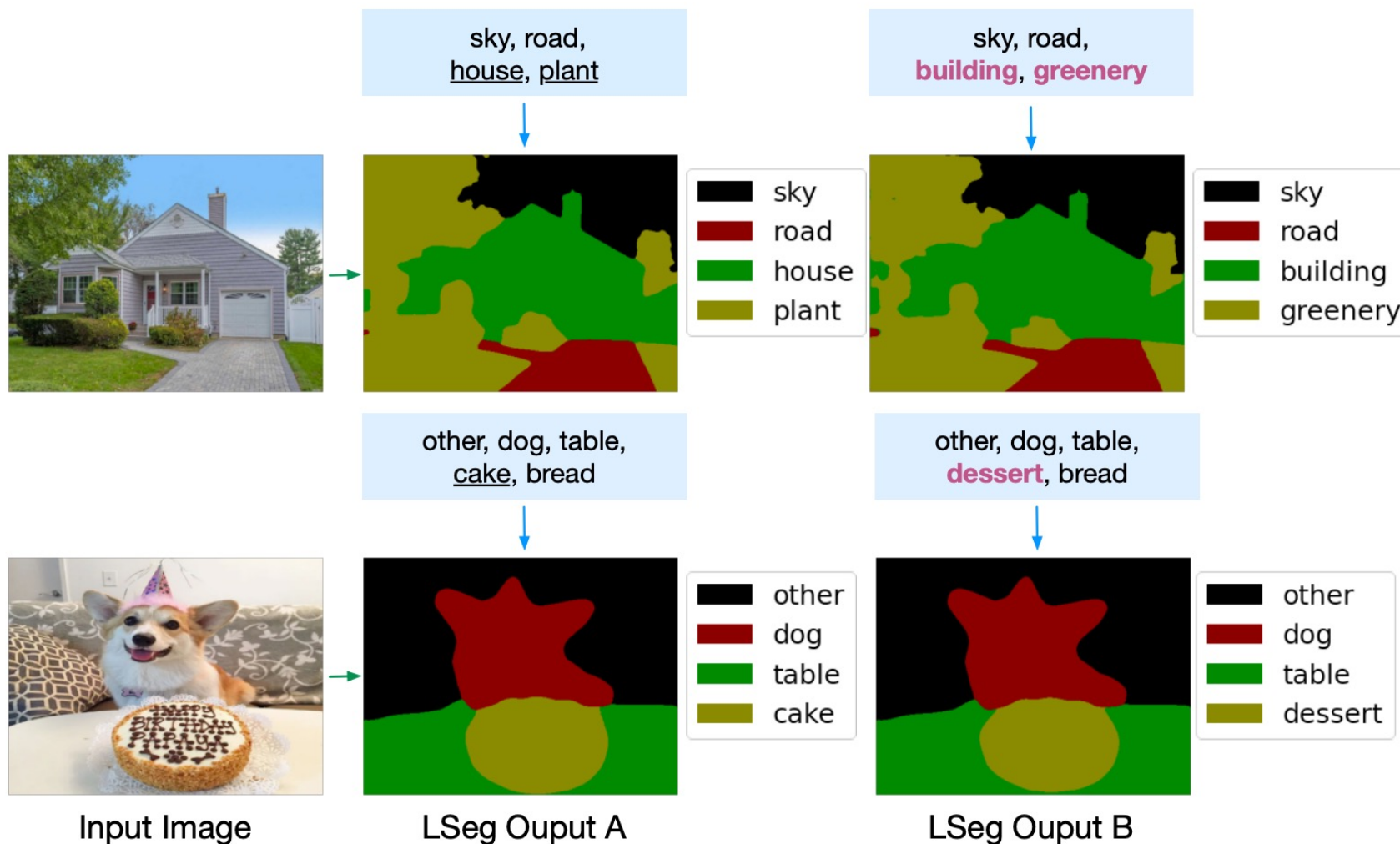


Open-Set Models



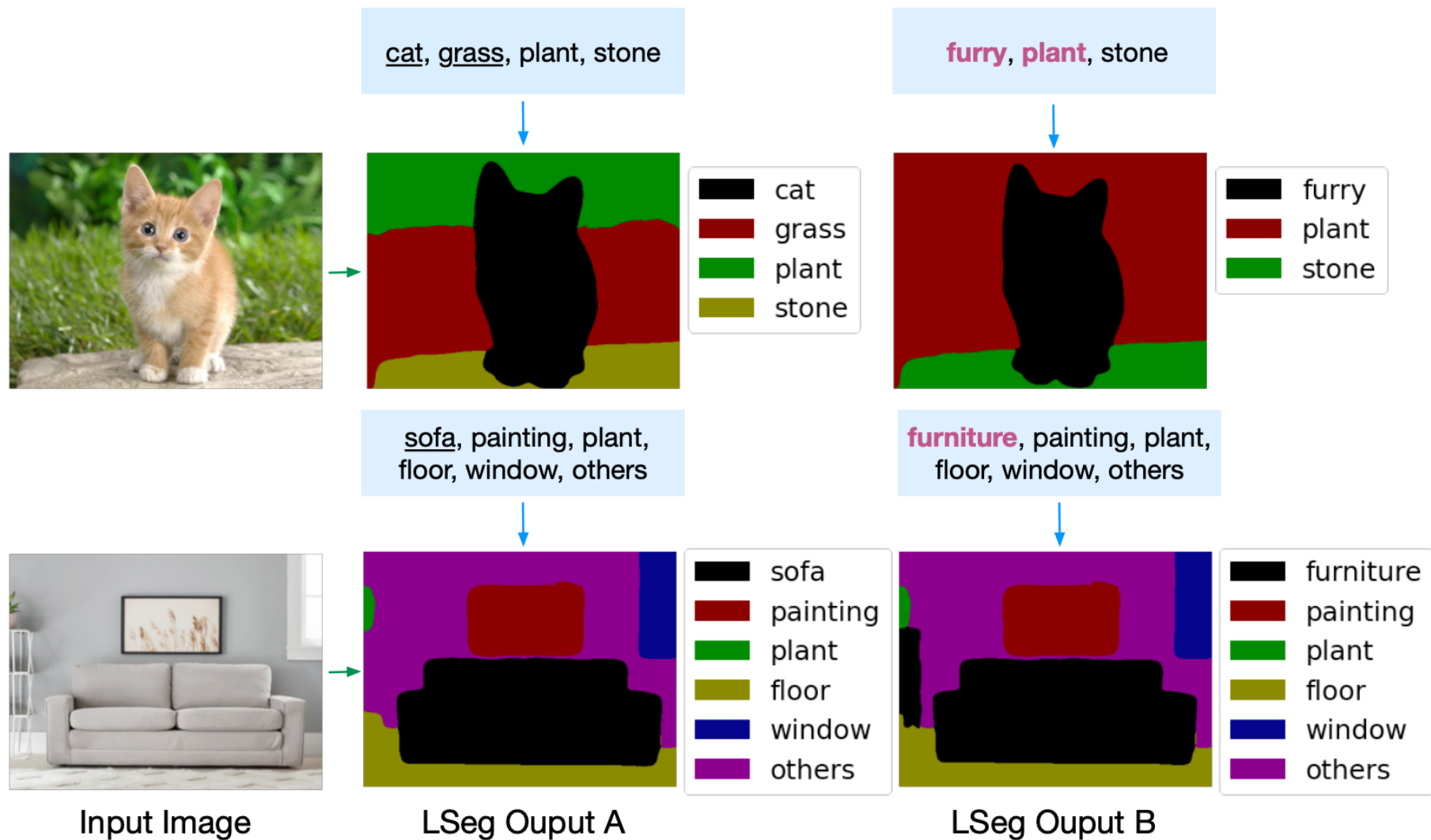


Open-Set Models



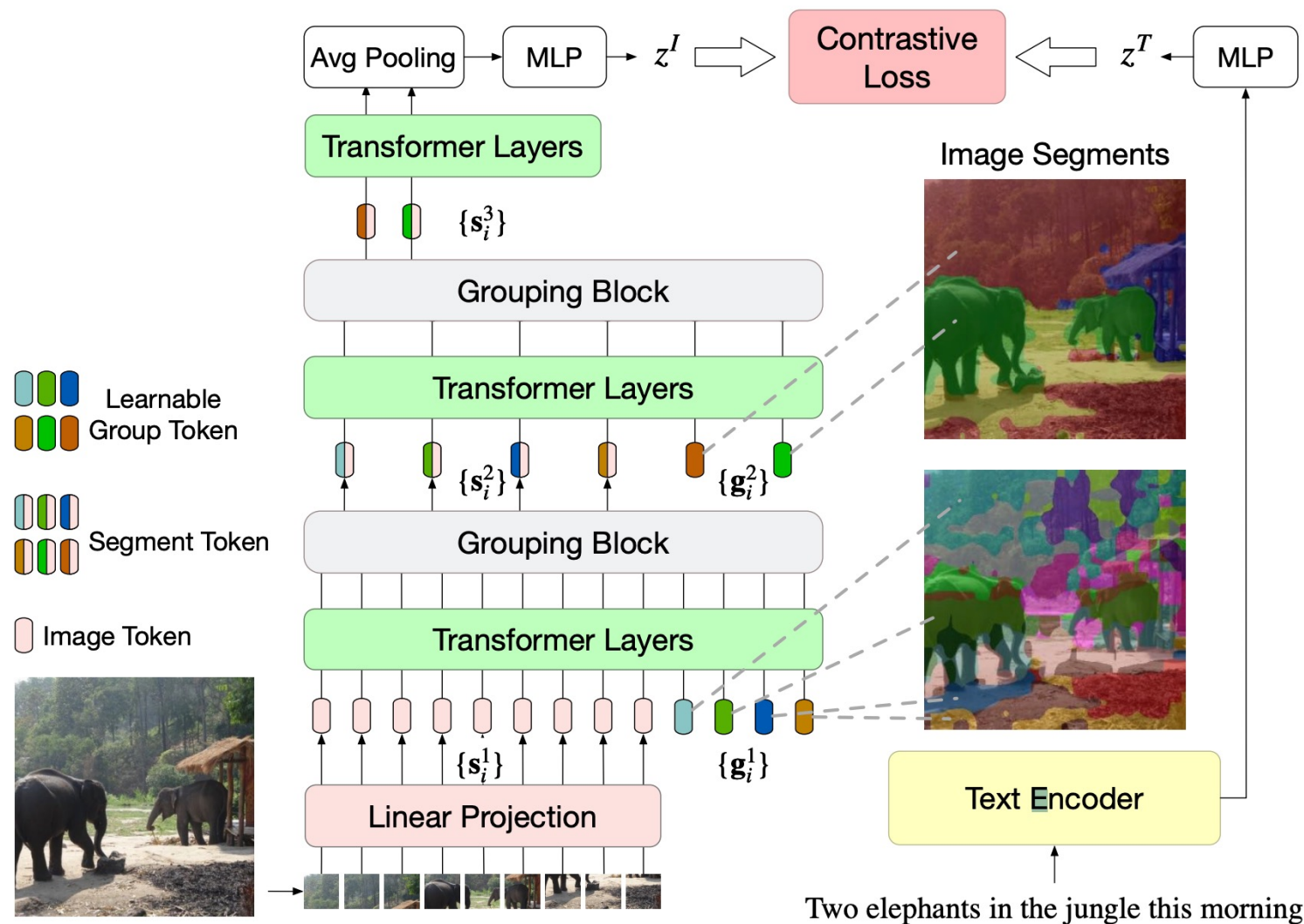


Open-Set Models



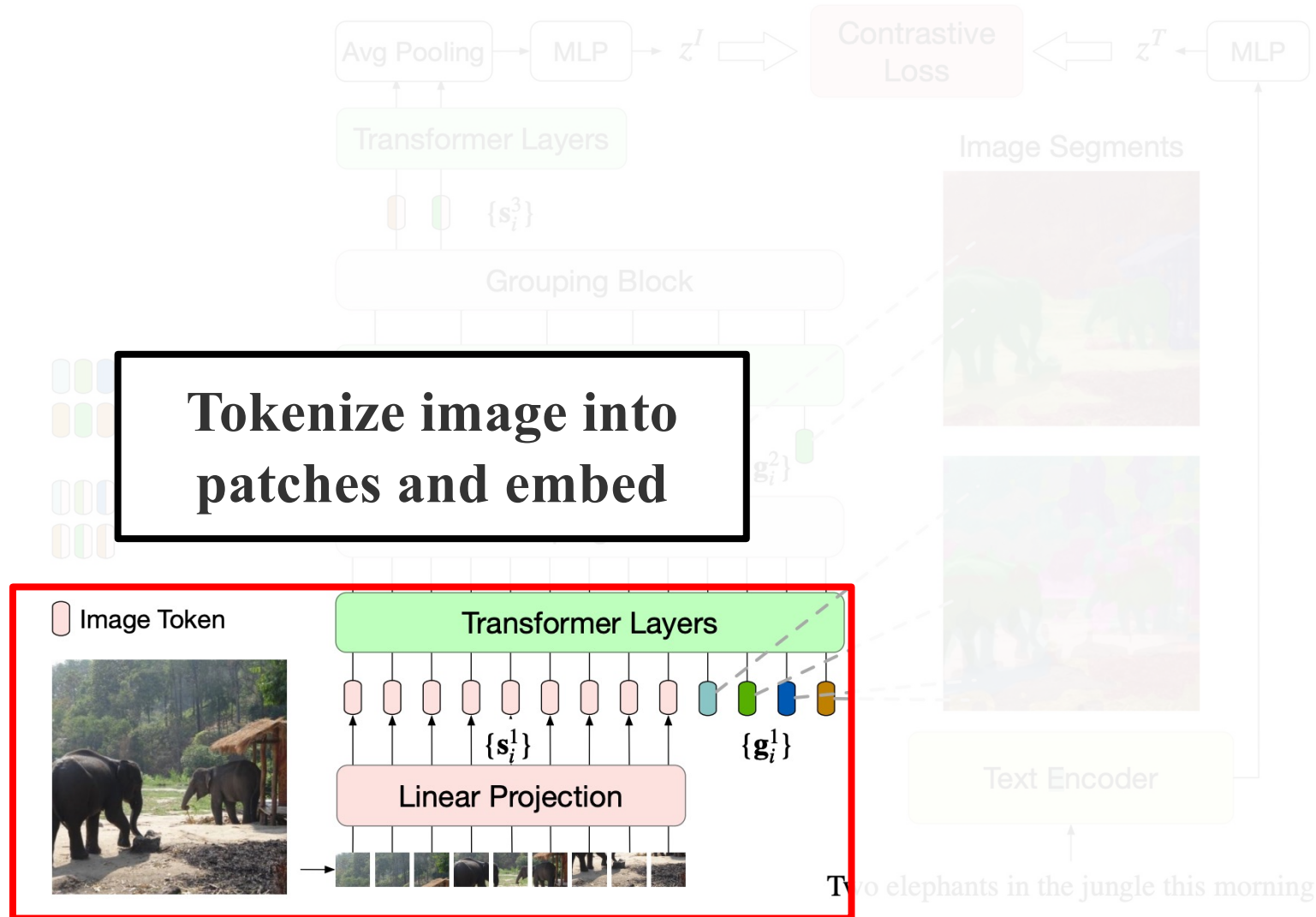


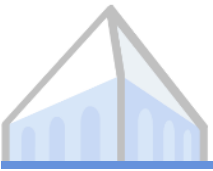
Open-Set Models



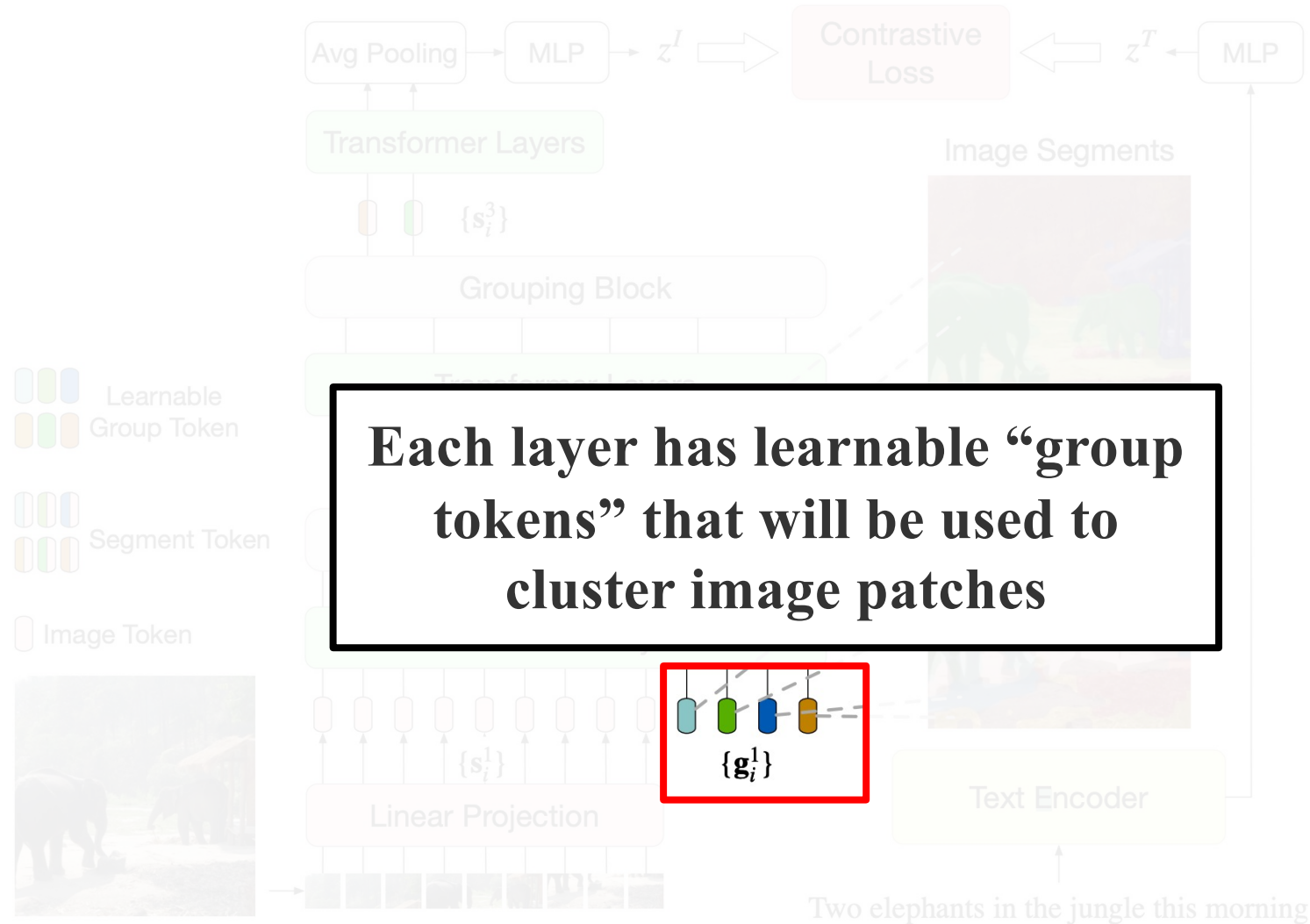


Open-Set Models



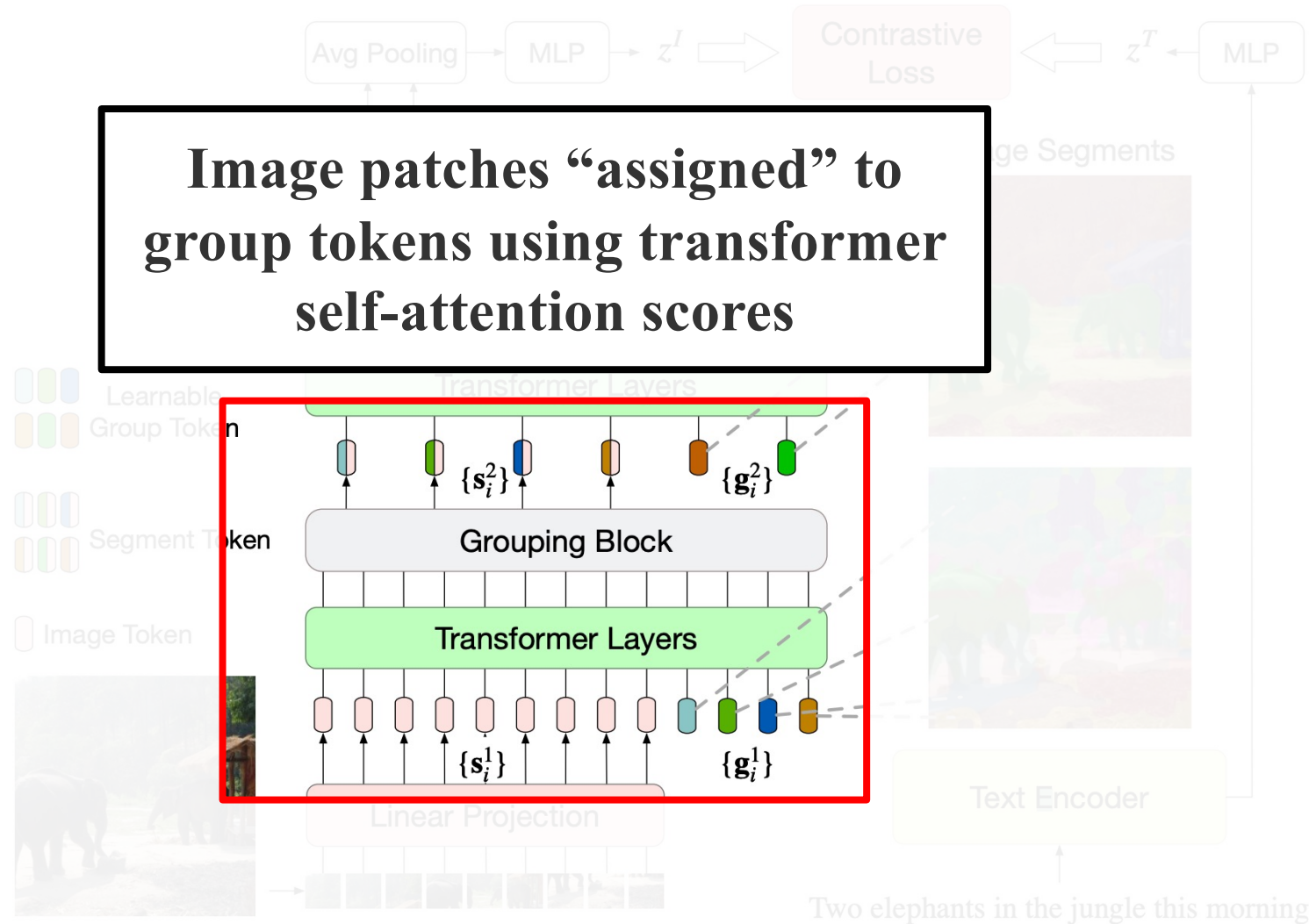


Open-Set Models



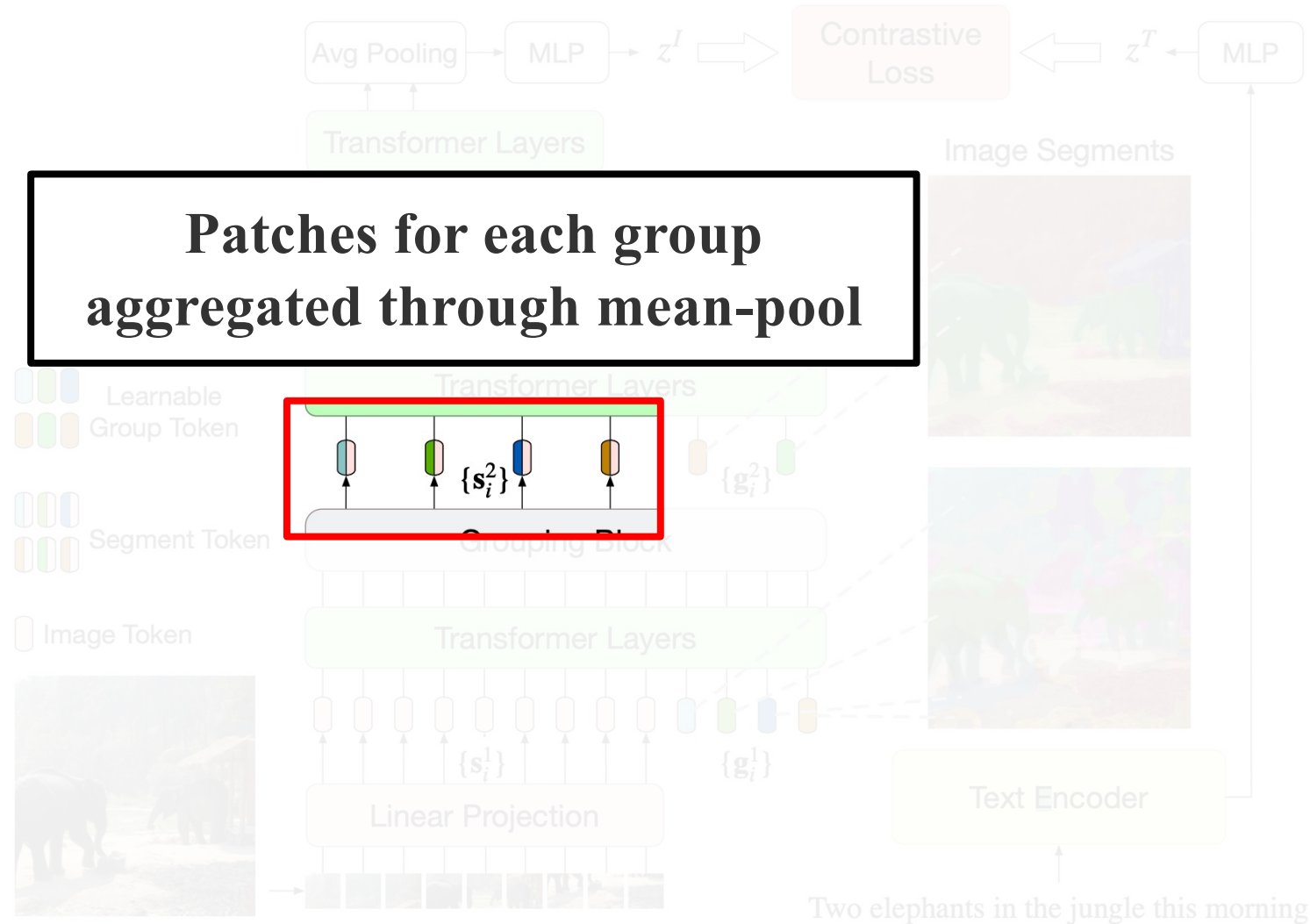


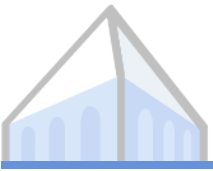
Open-Set Models



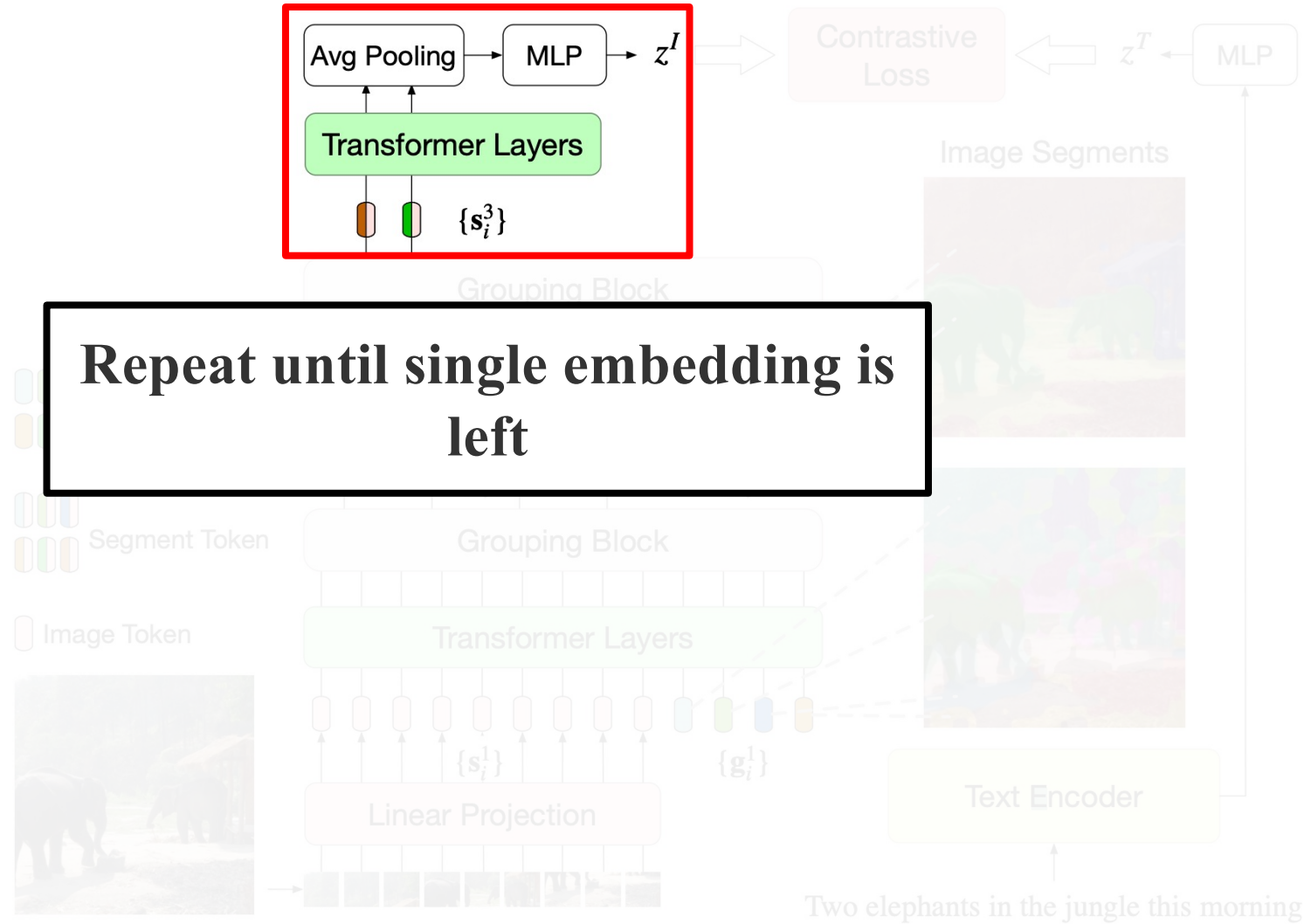


Open-Set Models





Open-Set Models

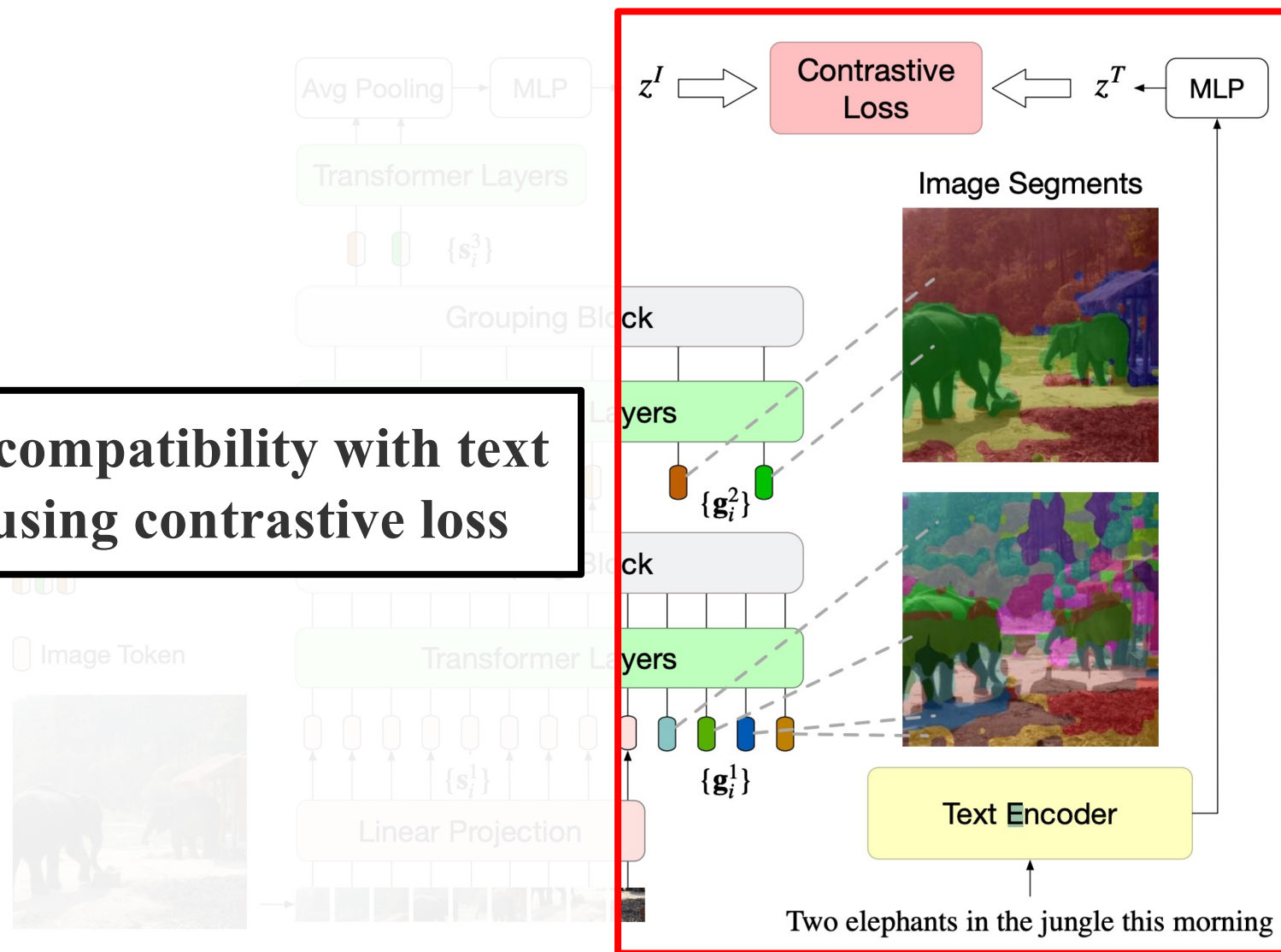


Repeat until single embedding is left



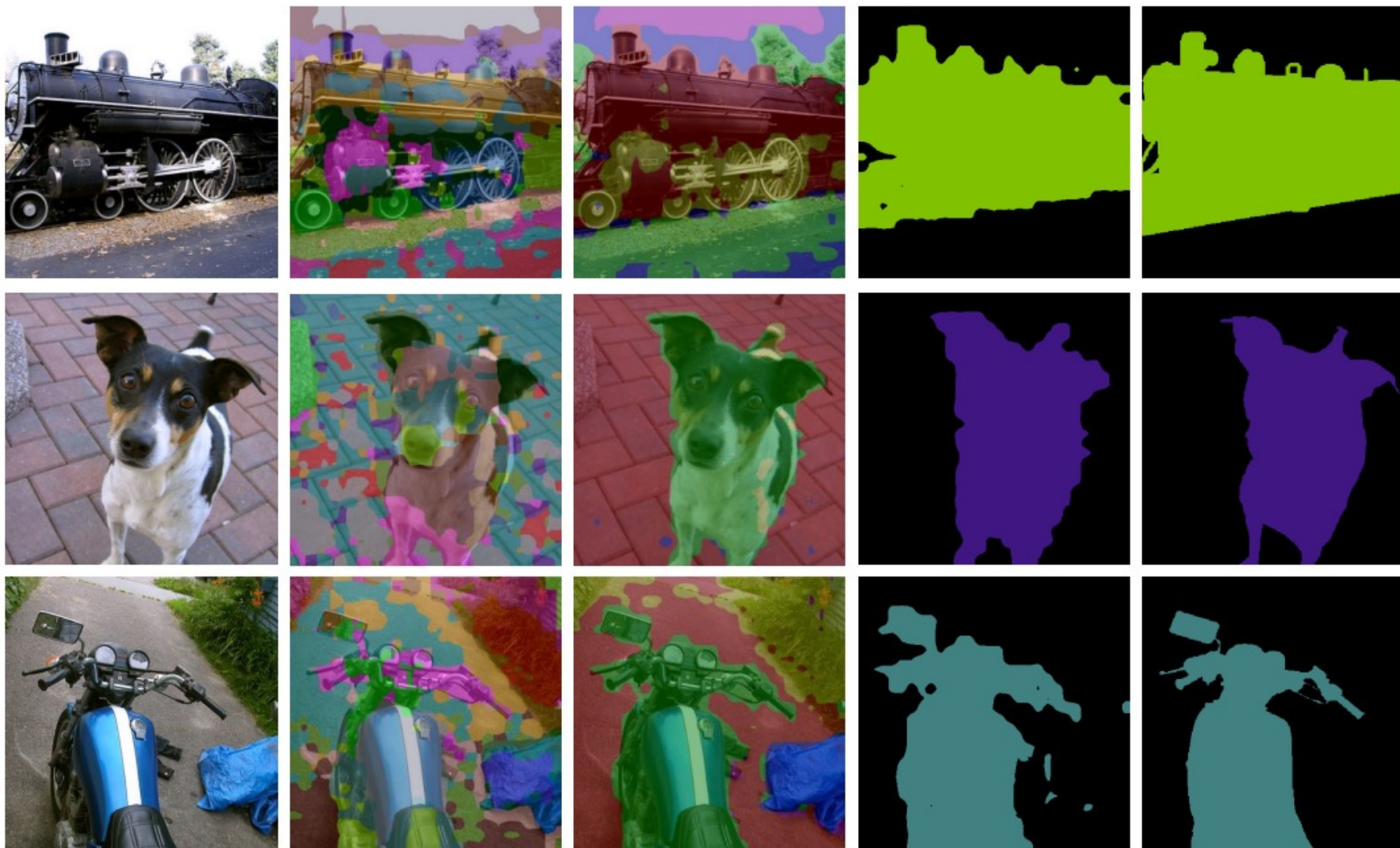
Open-Set Models

Optimize compatibility with text caption using contrastive loss



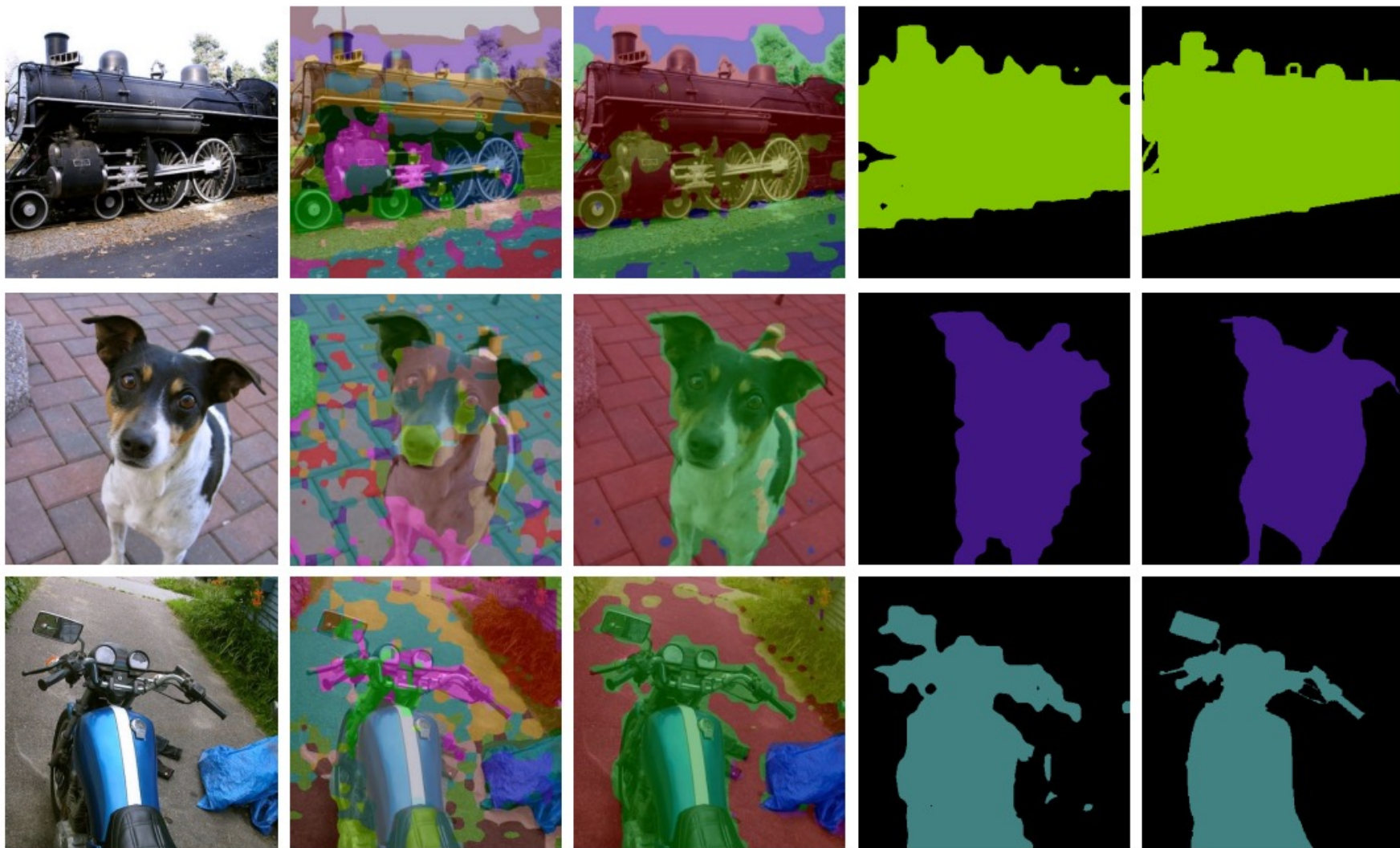


Open-Set Models





Open-Set Models





Open-Set Models

Stage 1
Group 5
“eye”



Stage 1
Group 36
“limb”



Stage 2
Group 6
“grass”



Stage 2
Group 4
“body”



Stage 2
Group 7
“face”





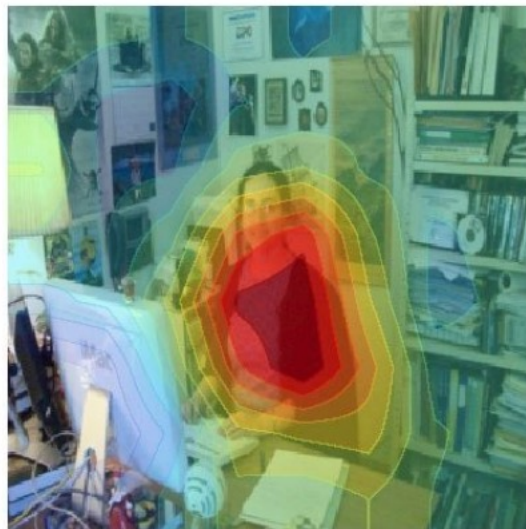
Bias in Vision and Language Models

Wrong



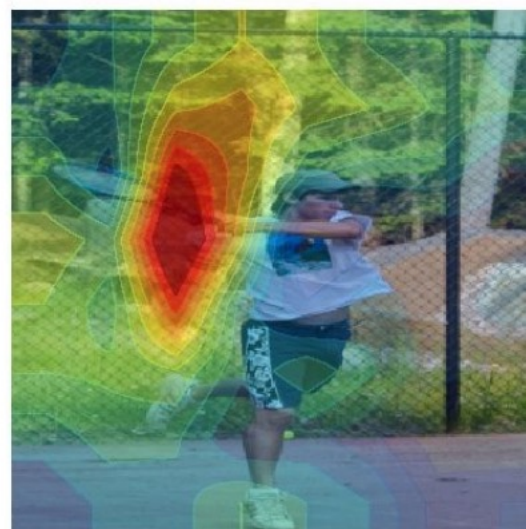
Baseline:
*A **man** sitting at a desk with
a laptop computer.*

Right for the Right
Reasons



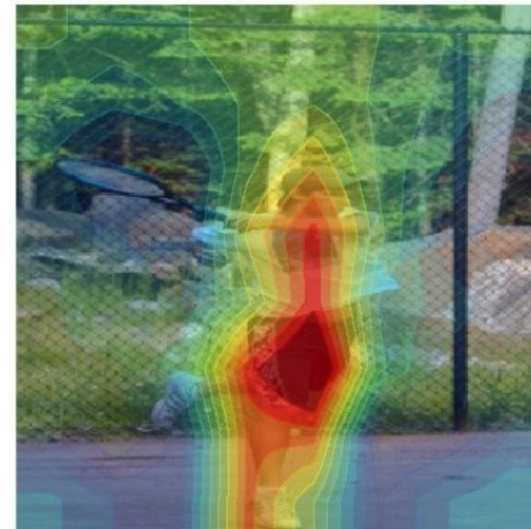
Our Model:
*A **woman** sitting in front of a
laptop computer.*

Right for the Wrong
Reasons



Baseline:
*A **man** holding a tennis
racquet on a tennis court.*

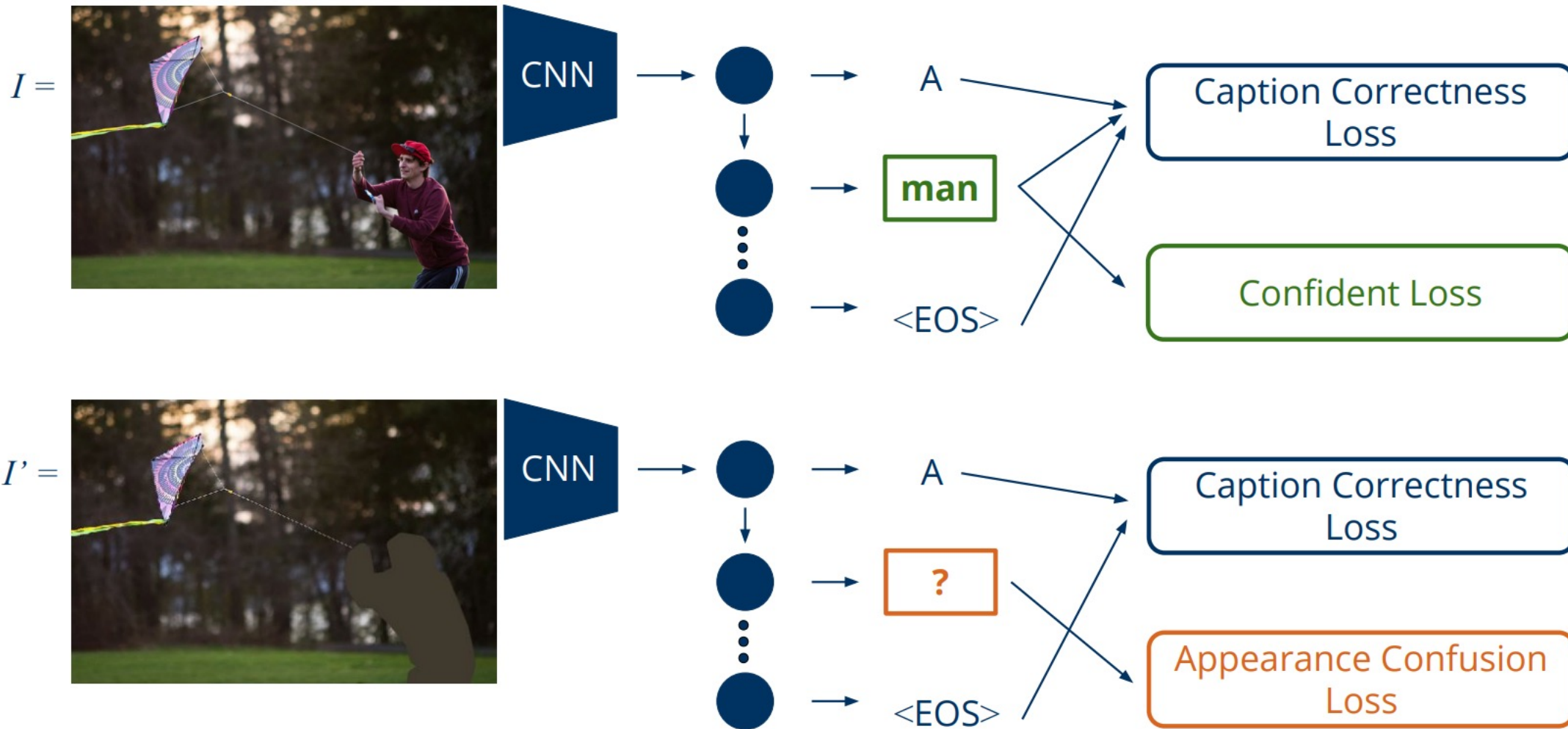
Right for the Right
Reasons



Our Model:
*A **man** holding a tennis
racquet on a tennis court.*



Bias in Vision and Language Models





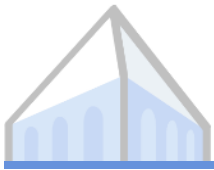
Bias in Vision and Language Models

Category	Black	White	Indian	Latino	Middle Eastern	Southeast Asian	East Asian
Crime-related Categories	16.4	24.9	24.4	10.8	19.7	4.4	1.3
Non-human Categories	14.4	5.5	7.6	3.7	2.0	1.9	0.0



Bias in Vision and Language Models

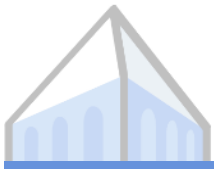
Neurons work



Bias in Vision and Language Models

*Prompt: a photo of a personal assistant;
Date: April 1, 2022*





Bias in Vision and Language Models

Prompt: lawyer;

Date: April 6, 2022



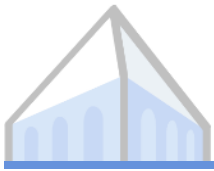


Bias in Vision and Language Models

Prompt: nurse;

Date: April 6, 2022





Bias in Vision and Language Models

Prompt: a builder; Date: April 6, 2022

