

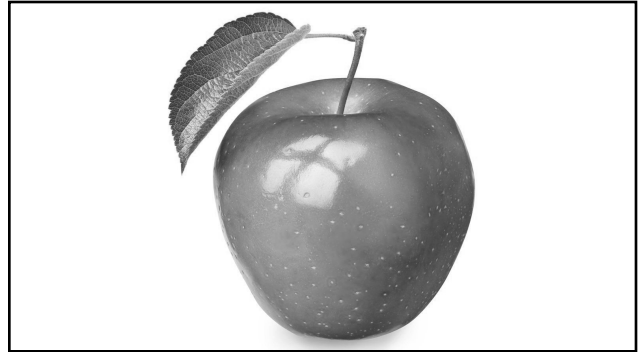
Vision and Language



Rodolfo (Rudy) Corona

with thanks to Daniel Fried  
CS 288, 4/12/2022

1

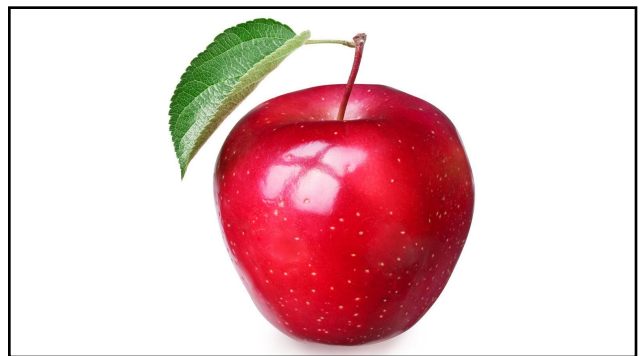


2

The colors of the visible light spectrum <sup>[1]</sup>		
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

Source: "Color" in Wikipedia

3



4

“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 ...”


5

### 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*.

6


### Grounding

- (Some) possible things to ground into:

7

### 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...



8

### 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



9

### 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



10

### 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




11

### 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]

12



## 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


13



## CS294-43: VISION AND LANGUAGE AI SEMINAR

---

14



## A Gallery of Tasks

15



## 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.

Microsoft COCO Captions: Chen et al. 2015


16

### 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




VQA 2.0: Goyal et al. 2017

17


### Object Detection (2D)



(a) Query: "street lamp"



(b) Query: "major league logo"

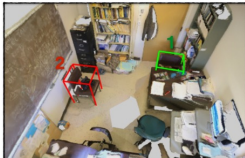


(c) Query: "zebras on savanna"


MDETR: Karnath et al. 2021

18

### 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."


ReferIt3D: Achlioptas et al. 2020

19

### Conditional Generation (2D)



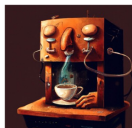
a classic portrait painting of Salvador Dalí with a colorful, half face




a dalmatian dog wearing a black beret and black turtleneck



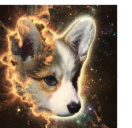
a close up of a hand holding a green sprout growing from a seed



an espresso machine that makes coffee from banana milk, vanilla



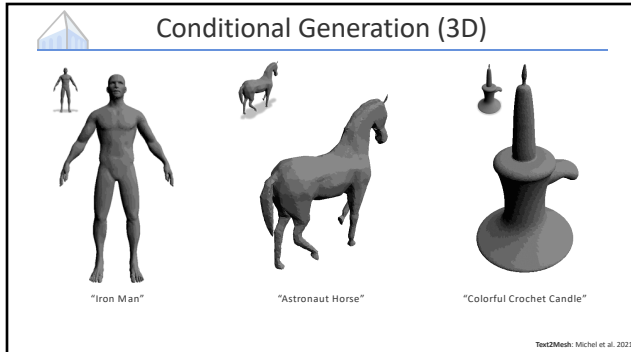
panda scientist wearing safety goggles, lab coat, and mask



a dog's head depicted as an explosion of energy

DALL-E 2: Ramesh et al. 2022

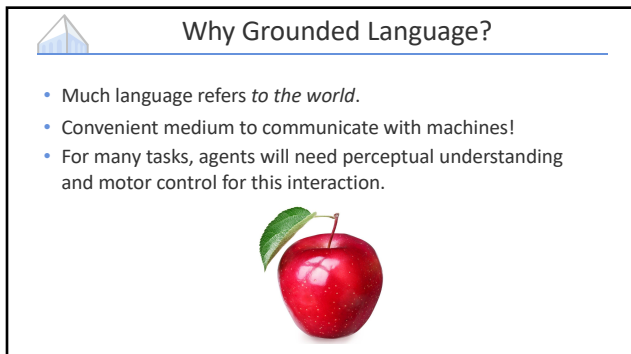
20



21



22



23



24



25

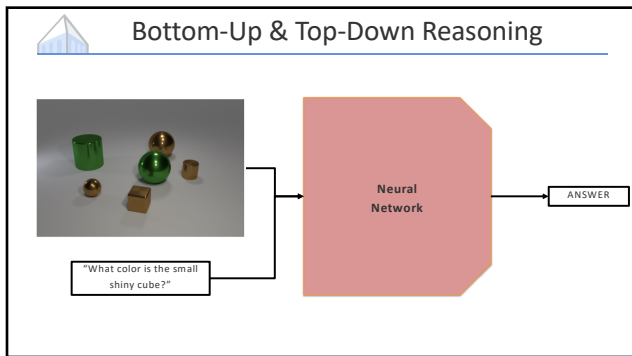
### Bottom-Up & Top-Down Reasoning

"What color is the small shiny cube?"

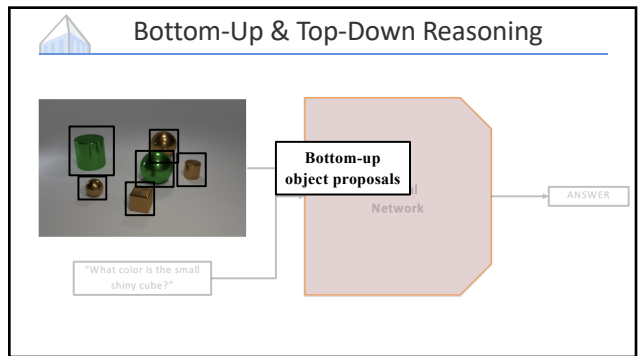
CLEVR: Johnson et al. 2016

A slide titled "Bottom-Up & Top-Down Reasoning" featuring a scene with several objects: a green cup, a green sphere, a brown cube, and a small shiny cube. Below the scene is a text box with the question "What color is the small shiny cube?". The slide includes a small blue pyramid icon in the top left and a small logo in the bottom right.

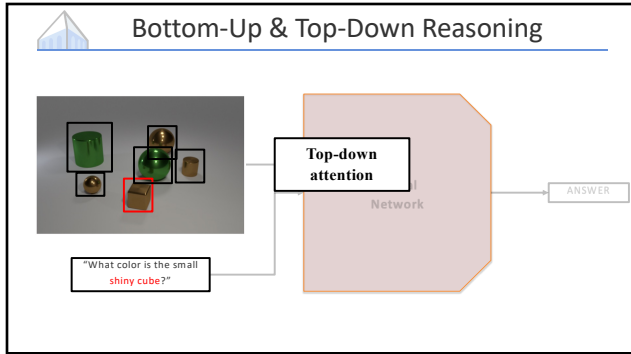
26



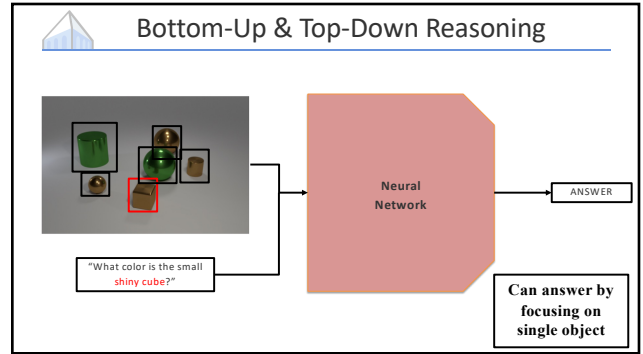
27



28



29



30

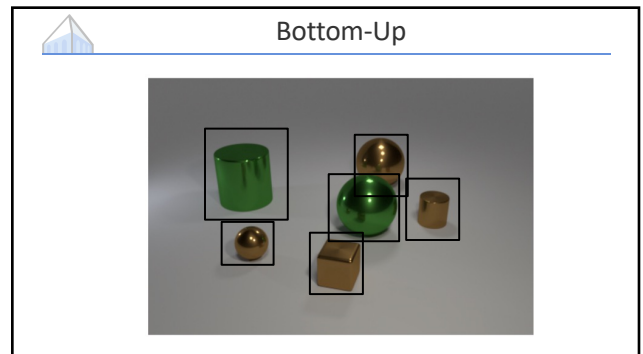
### Bottom-Up & Top-Down Reasoning

	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	<b>80.3</b>	<b>42.8</b>	<b>55.8</b>	<b>63.2</b>
Relative Improvement	3%	14%	8%	6%

**Provides inductive bias in both directions!**

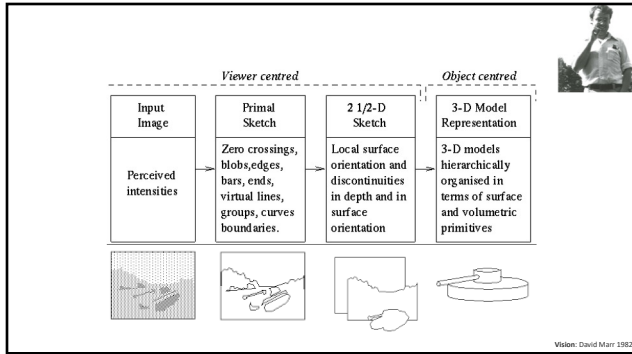
Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. Anderson et al. 2018

31

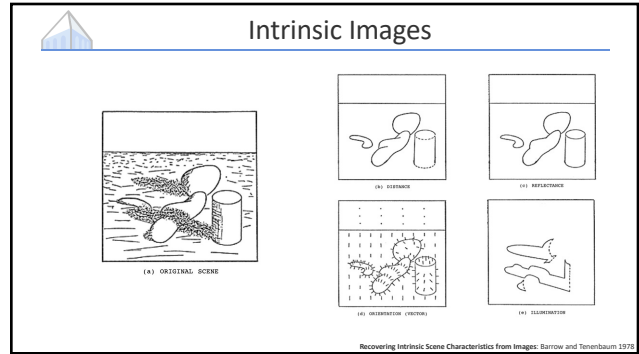


32





33



34

### "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

A Game-Theoretic Approach to Generating Spatial Descriptions: Gotland et al. 2010

35

### "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

**Task:** Describe target object unambiguously

A Game-Theoretic Approach to Generating Spatial Descriptions: Gotland et al. 2010

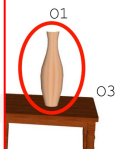
36

### “Solved” Perception

Question: Where is the object outlined in red?

Answer: The object outlined in red is O2

Relationships between objects known

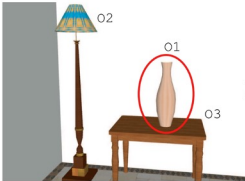


- left of
- right of
- above
- below
- in front of
- behind
- inside of
- on
- under
- across from

A Game-Theoretic Approach to Generating Spatial Descriptions: Golland et al. 2010

37

### “Solved” Perception



$w$

Right of O2  $\left\{ \begin{array}{l} O1 (p_L=0.5) \checkmark \\ O3 (p_L=0.5) \times \end{array} \right.$

$S(L)$ :

On top of O3  $\left\{ \begin{array}{l} O1 (p_L=1.0) \checkmark \end{array} \right.$

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


Problem reduced to pragmatic reasoning

A Game-Theoretic Approach to Generating Spatial Descriptions: Golland et al. 2010

38

### “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)]


Walk the Talk: MacMahon et al. 2006

39

### “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)]

Walk the Talk: MacMahon et al. 2006

40

### “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)]

Walk the Talk: MacMahon et al. 2006

41

### “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)]

Walk the Talk: MacMahon et al. 2006

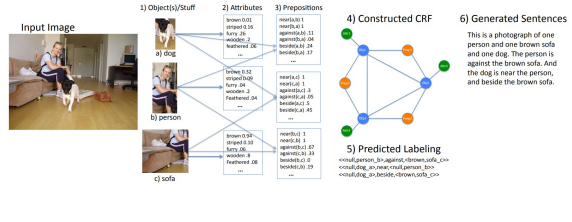
42

### “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.

43

### Intermediate Representations



Input Image

1) Object(s)/Stuff

2) Attributes

3) Prepositions

4) Constructed CRF

5) Predicted Labeling

6) Generated Sentences

BabyTalk: Kulkarni et al. 2013

44



### Intermediate Representations

Language model never sees pixels!

**1) Object(s)/Stuff**

- a) dog
- b) person
- c) sofa

**2) Attributes**

- brown: 0.01
- round: 0.11
- ... (for dog)
- ... (for person)
- ... (for sofa)

**3) Prepositions**

- near: 0.11
- on: 0.11
- ... (for dog)
- ... (for person)
- ... (for sofa)

**4) Constructed CRF**

**5) Predicted Labeling**

```
<<tag:person, p> agent: chow, s, 0.00
<<tag:dog, p> agent: mal, person, 0.00
<<tag: sofa, p> agent: brown, sofa, 0.00
```

**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.

BabyTalk: Kulkarni et al. 2013

49

### 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 blue sky.

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.

BabyTalk: Kulkarni et al. 2013

50

### Intermediate Representations

**The State Machine**

**alphabet (concepts)**

- bowl**: Color: brown (0.92), Material: wood (0.0)
- apple**: Color: red (0.95), Shape: round (0.87)
- girl**: Mood: happy (0.78), Posture: sitting (0.82)

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

**properties**: yellow, banana, grapes, bowl, table, right, coffee maker, man, behind, apple, inside, right, looking, girl, boy, smiling, states, transitions

**disentangled representation**: s<sub>1</sub> s<sub>2</sub> s<sub>3</sub> s<sub>4</sub> s<sub>5</sub> s<sub>6</sub> s<sub>7</sub> s<sub>8</sub> s<sub>9</sub> s<sub>10</sub> s<sub>11</sub> s<sub>12</sub> s<sub>13</sub> s<sub>14</sub> s<sub>15</sub> s<sub>16</sub> s<sub>17</sub> s<sub>18</sub> s<sub>19</sub> s<sub>20</sub> s<sub>21</sub> s<sub>22</sub> s<sub>23</sub> s<sub>24</sub> s<sub>25</sub> s<sub>26</sub> s<sub>27</sub> s<sub>28</sub> s<sub>29</sub> s<sub>30</sub> s<sub>31</sub> s<sub>32</sub> s<sub>33</sub> s<sub>34</sub> s<sub>35</sub> s<sub>36</sub> s<sub>37</sub> s<sub>38</sub> s<sub>39</sub> s<sub>40</sub> s<sub>41</sub> s<sub>42</sub> s<sub>43</sub> s<sub>44</sub> s<sub>45</sub> s<sub>46</sub> s<sub>47</sub> s<sub>48</sub> s<sub>49</sub> s<sub>50</sub> s<sub>51</sub> s<sub>52</sub> s<sub>53</sub> s<sub>54</sub> s<sub>55</sub> s<sub>56</sub> s<sub>57</sub> s<sub>58</sub> s<sub>59</sub> s<sub>60</sub> s<sub>61</sub> s<sub>62</sub> s<sub>63</sub> s<sub>64</sub> s<sub>65</sub> s<sub>66</sub> s<sub>67</sub> s<sub>68</sub> s<sub>69</sub> s<sub>70</sub> s<sub>71</sub> s<sub>72</sub> s<sub>73</sub> s<sub>74</sub> s<sub>75</sub> s<sub>76</sub> s<sub>77</sub> s<sub>78</sub> s<sub>79</sub> s<sub>80</sub> s<sub>81</sub> s<sub>82</sub> s<sub>83</sub> s<sub>84</sub> s<sub>85</sub> s<sub>86</sub> s<sub>87</sub> s<sub>88</sub> s<sub>89</sub> s<sub>90</sub> s<sub>91</sub> s<sub>92</sub> s<sub>93</sub> s<sub>94</sub> s<sub>95</sub> s<sub>96</sub> s<sub>97</sub> s<sub>98</sub> s<sub>99</sub> s<sub>100</sub>

Learning by Abstraction: The Neural State Machine: Hudson and Manning 2019

51

### Intermediate Representations

**The State Machine**

**alphabet (concepts)**

- bowl**: Color: brown (0.92), Material: wood (0.0)
- apple**: Color: red (0.95), Shape: round (0.87)
- girl**: Mood: happy (0.78), Posture: sitting (0.82)

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

**properties**: yellow, banana, grapes, bowl, table, right, coffee maker, man, behind, apple, inside, right, looking, girl, boy, smiling, states, transitions

**disentangled representation**: s<sub>1</sub> s<sub>2</sub> s<sub>3</sub> s<sub>4</sub> s<sub>5</sub> s<sub>6</sub> s<sub>7</sub> s<sub>8</sub> s<sub>9</sub> s<sub>10</sub> s<sub>11</sub> s<sub>12</sub> s<sub>13</sub> s<sub>14</sub> s<sub>15</sub> s<sub>16</sub> s<sub>17</sub> s<sub>18</sub> s<sub>19</sub> s<sub>20</sub> s<sub>21</sub> s<sub>22</sub> s<sub>23</sub> s<sub>24</sub> s<sub>25</sub> s<sub>26</sub> s<sub>27</sub> s<sub>28</sub> s<sub>29</sub> s<sub>30</sub> s<sub>31</sub> s<sub>32</sub> s<sub>33</sub> s<sub>34</sub> s<sub>35</sub> s<sub>36</sub> s<sub>37</sub> s<sub>38</sub> s<sub>39</sub> s<sub>40</sub> s<sub>41</sub> s<sub>42</sub> s<sub>43</sub> s<sub>44</sub> s<sub>45</sub> s<sub>46</sub> s<sub>47</sub> s<sub>48</sub> s<sub>49</sub> s<sub>50</sub> s<sub>51</sub> s<sub>52</sub> s<sub>53</sub> s<sub>54</sub> s<sub>55</sub> s<sub>56</sub> s<sub>57</sub> s<sub>58</sub> s<sub>59</sub> s<sub>60</sub> s<sub>61</sub> s<sub>62</sub> s<sub>63</sub> s<sub>64</sub> s<sub>65</sub> s<sub>66</sub> s<sub>67</sub> s<sub>68</sub> s<sub>69</sub> s<sub>70</sub> s<sub>71</sub> s<sub>72</sub> s<sub>73</sub> s<sub>74</sub> s<sub>75</sub> s<sub>76</sub> s<sub>77</sub> s<sub>78</sub> s<sub>79</sub> s<sub>80</sub> s<sub>81</sub> s<sub>82</sub> s<sub>83</sub> s<sub>84</sub> s<sub>85</sub> s<sub>86</sub> s<sub>87</sub> s<sub>88</sub> s<sub>89</sub> s<sub>90</sub> s<sub>91</sub> s<sub>92</sub> s<sub>93</sub> s<sub>94</sub> s<sub>95</sub> s<sub>96</sub> s<sub>97</sub> s<sub>98</sub> s<sub>99</sub> s<sub>100</sub>

**Generate scene graph from image**

Learning by Abstraction: The Neural State Machine: Hudson and Manning 2019

52

### Intermediate Representations

**Graph vocabulary predefined**

**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*

Learning by Abstraction: The Neural State Machine: Hudson and Manning 2019

53

### Intermediate Representations

**Transform question into program traversing graph for answer**

*instructions*     *disentangled representation*

Learning by Abstraction: The Neural State Machine: Hudson and Manning 2019

54

### Intermediate Representations

**Answer by executing program in state machine**

*instructions*     *disentangled representation*

Learning by Abstraction: The Neural State Machine: Hudson and Manning 2019

55

### Intermediate Representations

**Allows language reasoning to occur solely within abstract structure**

*instructions*     *disentangled representation*

Learning by Abstraction: The Neural State Machine: Hudson and Manning 2019

56

### 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
<b>NSM</b>	<b>40.24</b>	<b>55.72</b>

Learning by Abstraction: The Neural State Machine: Hudson and Manning 2019

57

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:

DALL-E 1: Ramesh et al. 2021

58

### 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

↕

DALL-E 1: Ramesh et al. 2021

59

### 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 Dord et al. 2017

DALL-E 1: Ramesh et al. 2021

60

### 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

DALL-E 1: Ramesh et al. 2021

61

### 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

DALL-E 1: Ramesh et al. 2021

62

### 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!

DALL-E 1: Ramesh et al. 2021

63

TEXT PROMPT: an x-ray of a capybara sitting in a forest

AI-GENERATED IMAGES

DALL-E 1: Ramesh et al. 2021

64



### 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”*

Image Understanding: John Torrioni 1987

65

### Anchoring to 3D

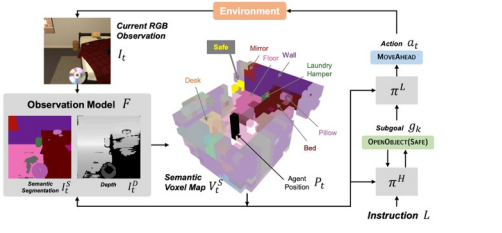


“Place a clean ladle on a counter”

ALFRED: Shridhar et al. 2020

66

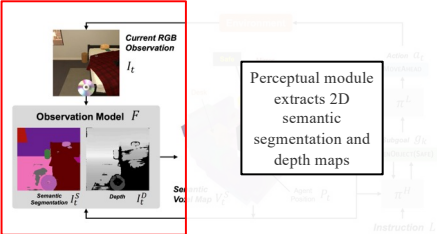
### Anchoring to 3D



A Persistent Spatial Semantic Representation for High-Level Natural Language Instruction Execution: Blukis et al. 2021

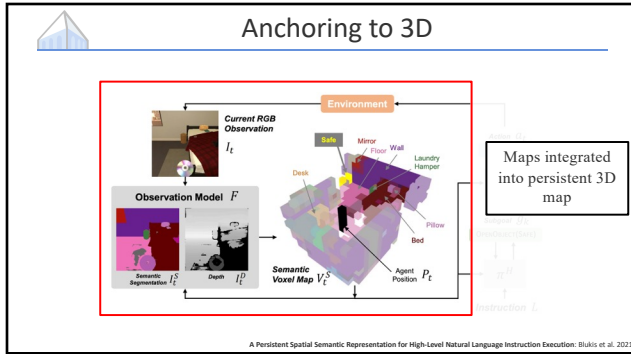
67

### Anchoring to 3D

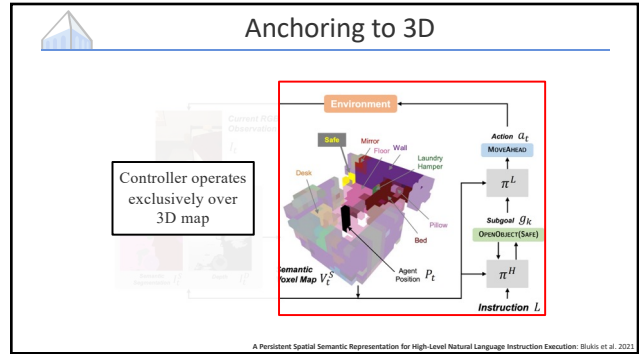


A Persistent Spatial Semantic Representation for High-Level Natural Language Instruction Execution: Blukis et al. 2021

68



69



70

### Anchoring to 3D

Method	Validation			
	Seen		Unseen	
	SR	GC	SR	GC
<b>HLSM</b>	29.6	38.8	18.3	<b>31.2</b>
+ 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	<b>18.9</b>	30.8

3D Map useful for improving performance

A Persistent Spatial Semantic Representation for High-Level Natural Language Instruction Execution. Blukis et al. 2021

71

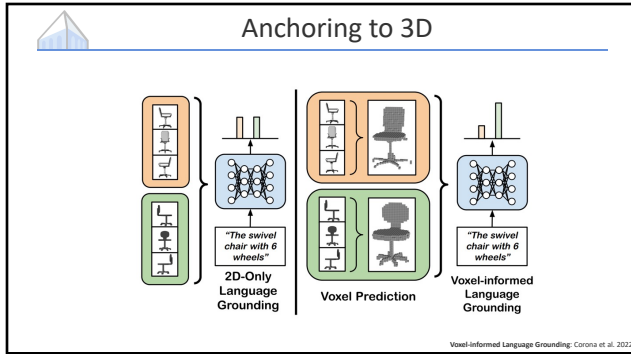
### Anchoring to 3D

Method	Validation			
	Seen		Unseen	
	SR	GC	SR	GC
HLSM	29.6	38.8	18.3	<b>31.2</b>
+ gt depth	29.6	40.5	20.1	33.7
<b>+ gt depth, gt seg.</b>	<b>40.7</b>	<b>50.4</b>	<b>40.2</b>	<b>52.2</b>
+ 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	<b>18.9</b>	30.8

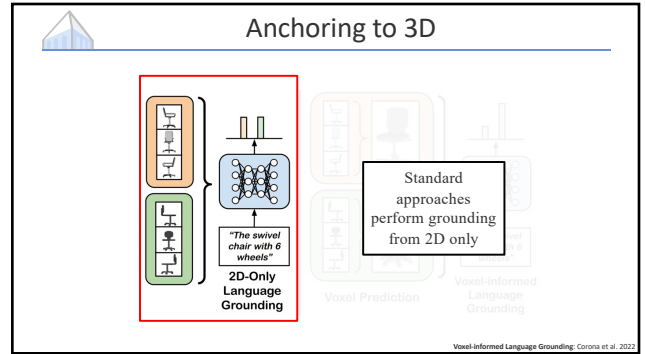
However, benefits held back by cascading errors

A Persistent Spatial Semantic Representation for High-Level Natural Language Instruction Execution. Blukis et al. 2021

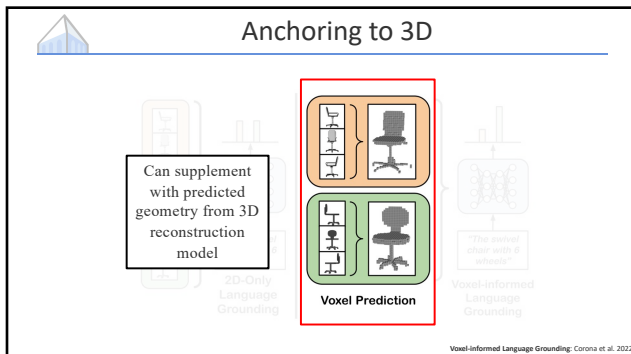
72



73



74



75

### Anchoring to 3D

Model	VALIDATION			TEST		
	Visual	Blind	All	Visual	Blind	All
ViLBERT	89.5	76.6	83.1	80.2	<b>73.0</b>	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)	-	-	-
<b>VLG (Ours)</b>	<b>91.2</b> (0.4)	<b>78.4</b> (0.7)	<b>84.9</b> (0.3)	<b>86.0</b>	71.7	<b>79.0</b>


Improves performance over 2D-only approaches

Voxel-informed Language Grounding: Corona et al. 2022

76

### 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.

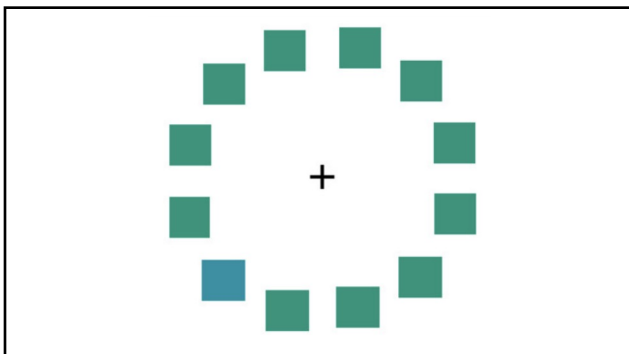


77

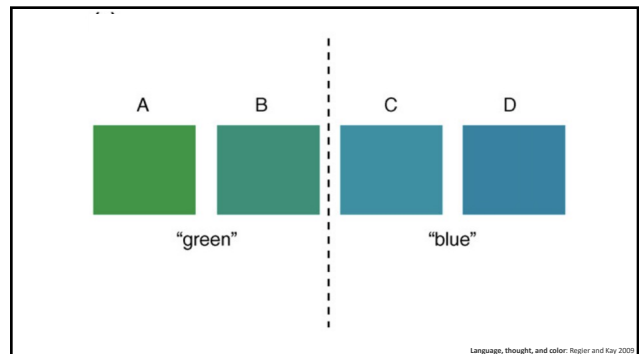
### Top-Down

“What color is the small shiny cube?”

78



79



80



81

**WordNet Search - 3.1**  
[- WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

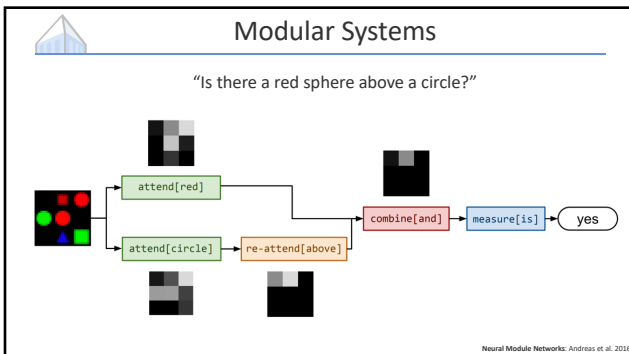
Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations  
 Display options for sense: (gloss) "an example sentence"

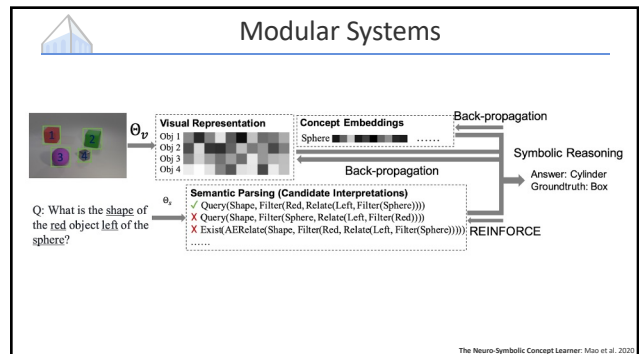
**Noun**

- 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)

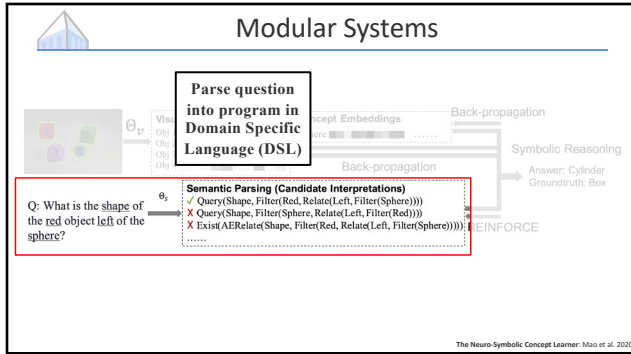
82



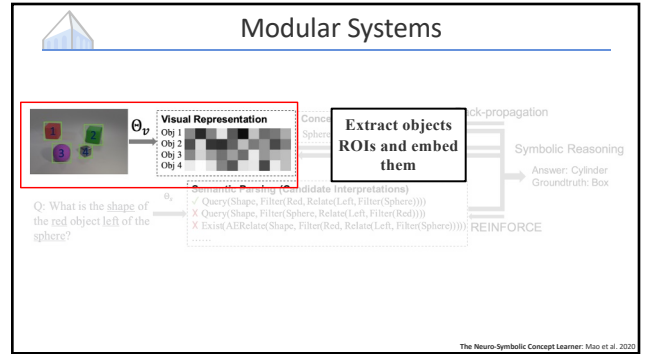
83



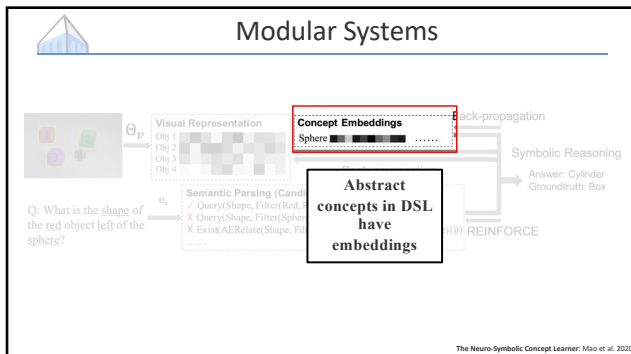
84



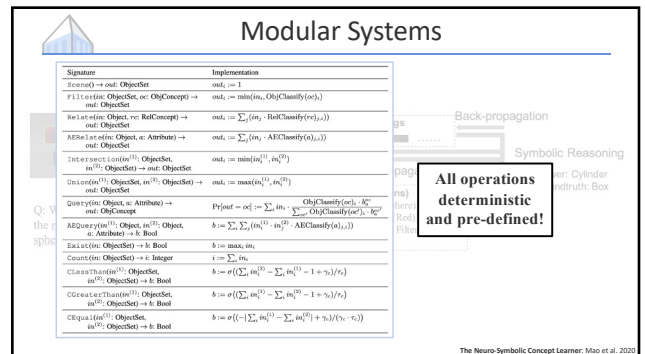
85



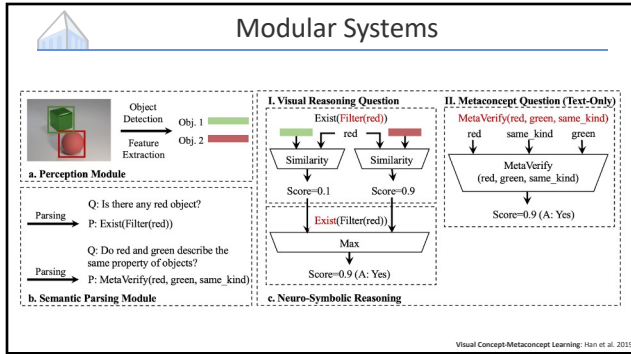
86



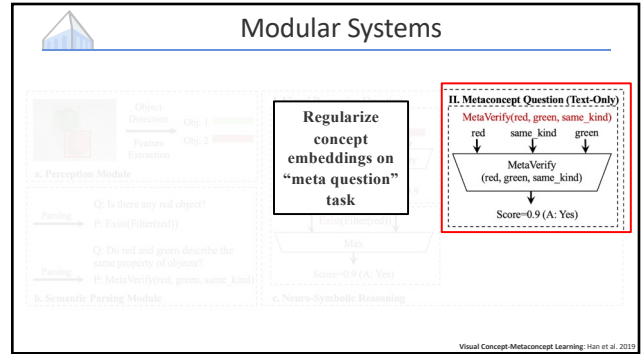
87



88



89



90

### Modular Systems

"block" == "square"

	GRU-CNN	MAC	NS-CL	VCML
<b>CLEVR</b>	50.0±0.0	68.7±3.8	80.2±3.1	<b>94.1±4.6</b>
<b>GQA</b>	50.0±0.0	49.5±0.2	49.3±0.6	<b>50.5±0.1</b>

Learning *synonyms* helps zero-shot generalization

Visual Concept-Metaconcept Learning: Han et al. 2019

91

### Modular Systems

■ == "purple" + "square"

	GRU-CNN	MAC	NS-CL	VCML
<b>CLEVR-200</b>	50.0±0.0	94.2±3.3	98.5±0.3	<b>98.9±0.2</b>
<b>CLEVR-20</b>	50.0±0.0	79.7±2.6	<b>95.7±0.0</b>	95.1±1.6

Learning *same kind* helps compositional generalization

Visual Concept-Metaconcept Learning: Han et al. 2019

92

### Modular Systems

"All the dogs are black."

Basic-NMN: find (dogs) → count (14) + count (16) → equal → False (57%)

Faithful-NMN: find (dogs) → filter (black) → count (9) + count (2) → equal → False (98%)

Obtaining Faithful Interpretations from Compositional Neural Networks: Subramanian et al. 2020

93

### Language as Signal for Abstractions

Learning with Latent Language: Andreas et al. 2017

94

### Language as Signal for Abstractions

Available at Training

a white shape is left of a yellow semicircle

Learning with Latent Language: Andreas et al. 2017

95

### Language as Signal for Abstractions

Meta (Snell et al., 2017): Support  $I_s$ , Query  $I_q$  → LSTM-Enc →  $g_s$  → LSTM-Dec → True

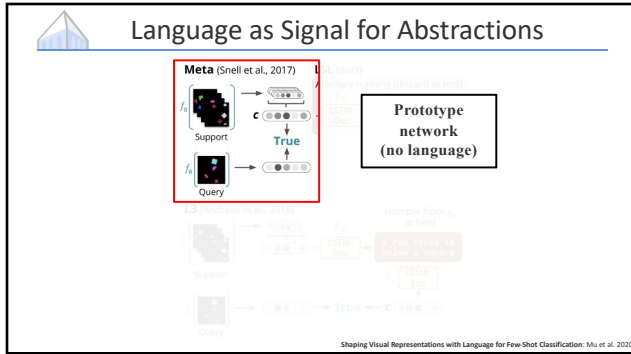
L3 (Andreas et al., 2018): Support  $I_s$ , Query  $I_q$  → LSTM-Enc →  $g_s$  → LSTM-Dec → True

LSL (ours): Auxiliary training (discard at test): Support  $I_s$  → LSTM-Enc →  $g_s$  → LSTM-Dec → a red cross is below a square

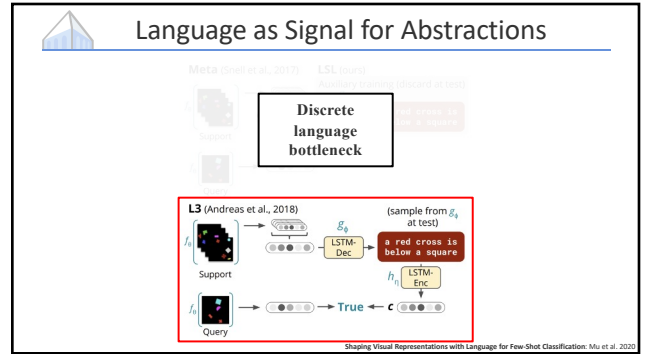
Shaping Visual Representations with Language for Few-Shot Classification: Mu et al. 2020

96

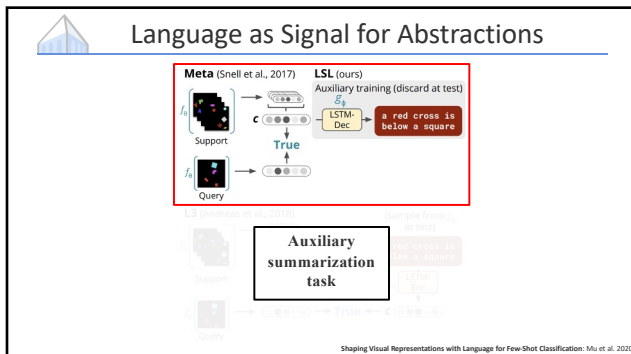




97



98



99

### Language as Signal for Abstractions

Test Set Accuracy

	ShapeWorld	Birds ( $D = 20$ )
Meta	60.59 ± 1.07	57.97 ± 0.96
L3	66.60 ± 1.18	53.96 ± 1.06
LSL	<b>67.29 ± 1.03</b>	<b>61.24 ± 0.96</b>

Shaping Visual Representations with Language for Few-Shot Classification: Mu et al. 2020

100


### 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](#)

101

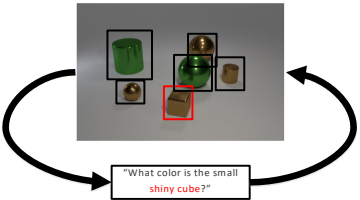
### Bottom-Up & Top-Down Reasoning



"What color is the small  
shiny cube?"


102

### Bottom-Up & Top-Down Reasoning



"What color is the small  
shiny cube?"

103



Galatea of the Spheres, Salvador Dali 1952

104

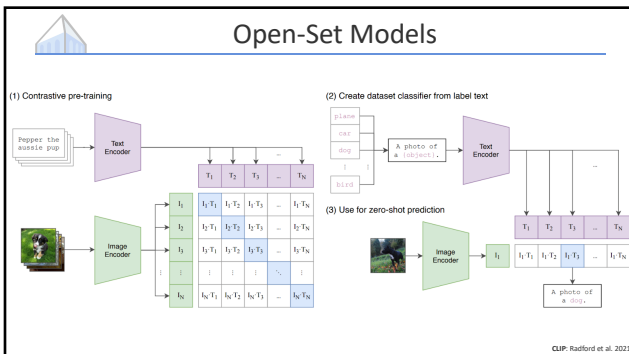
Extra Slides

105

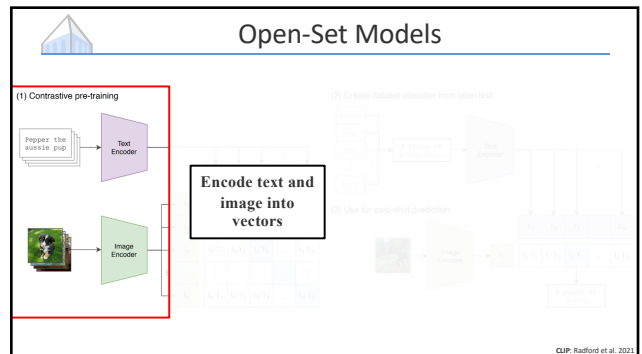
Open-Set Models

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

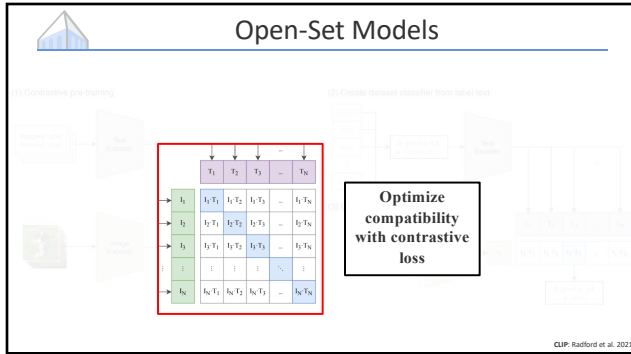
106



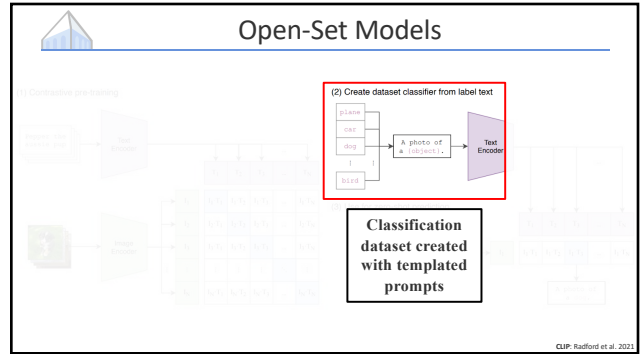
107



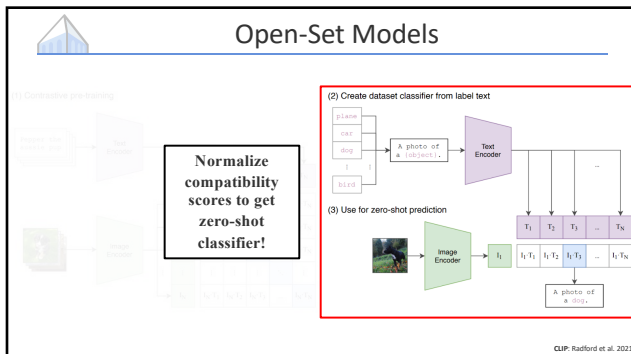
108



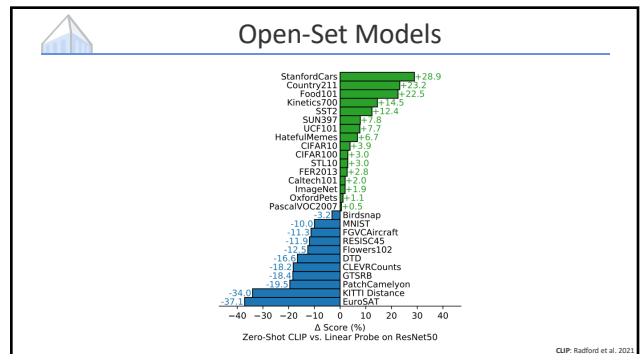
109



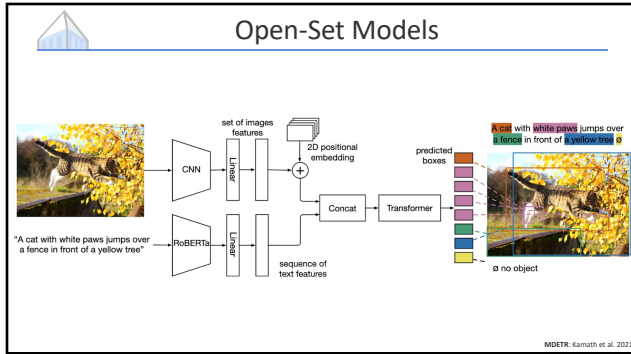
110



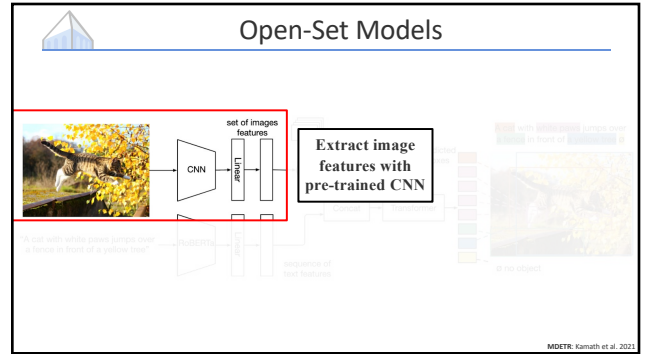
111



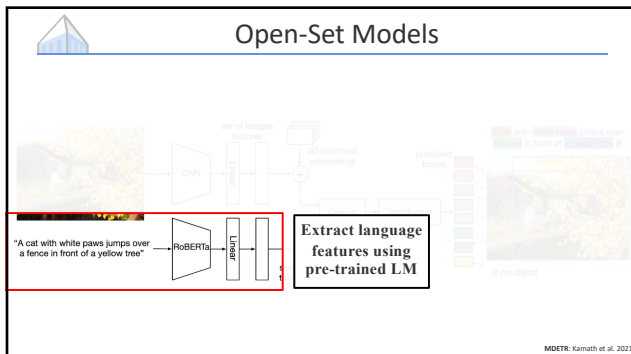
112



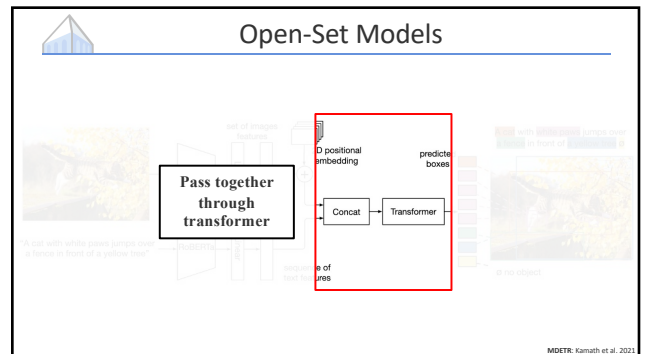
113



114



115



116

### Open-Set Models

Learned embedding "queries" tied to tokens in input text/image regions

MDETR: Kamath et al. 2021

117

### 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"

MDETR: Kamath et al. 2021

118

### Open-Set Models

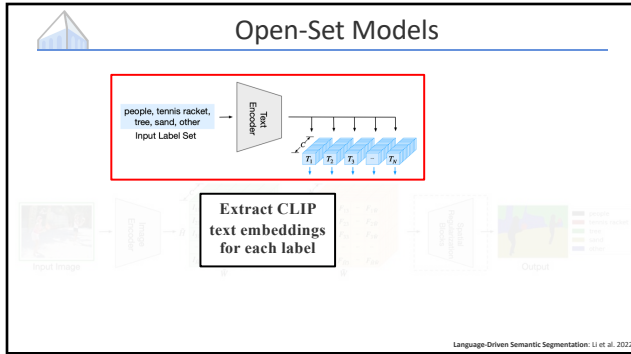
Language-Driven Semantic Segmentation: Li et al. 2022

119

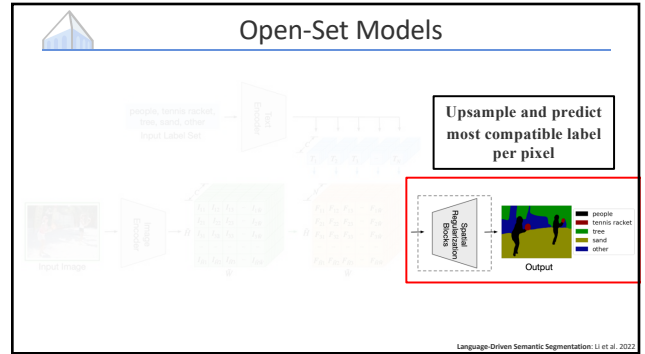
### Open-Set Models

Language-Driven Semantic Segmentation: Li et al. 2022

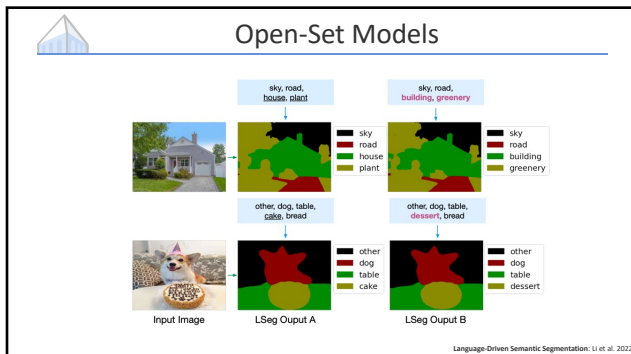
120



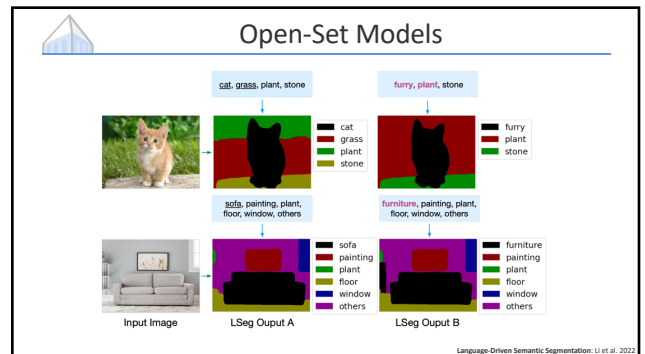
121



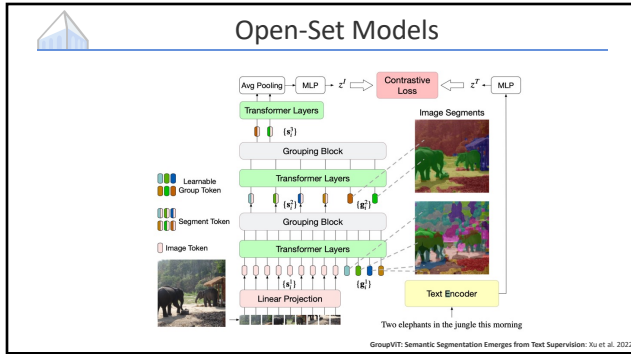
122



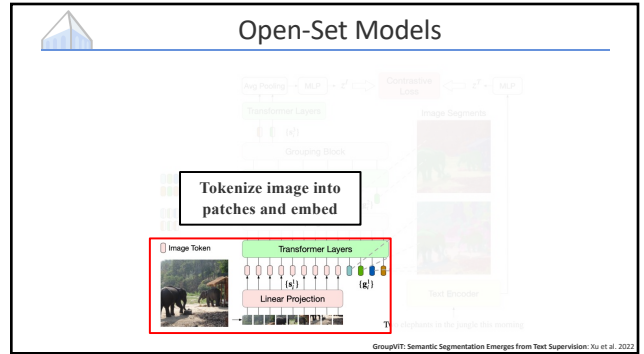
123



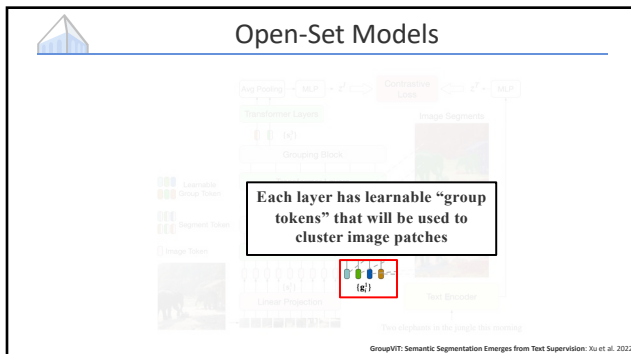
124



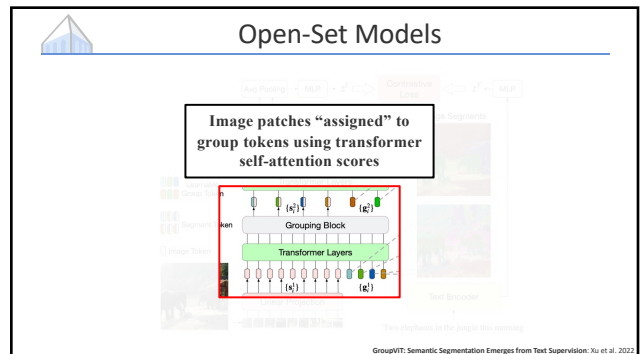
125



126



127



128



### Open-Set Models

**Patches for each group aggregated through mean-pool**

GroupViT: Semantic Segmentation Emerges from Text Supervision: Xu et al. 2022

129

### Open-Set Models

**Repeat until single embedding is left**

GroupViT: Semantic Segmentation Emerges from Text Supervision: Xu et al. 2022

130

### Open-Set Models

**Optimize compatibility with text caption using contrastive loss**

GroupViT: Semantic Segmentation Emerges from Text Supervision: Xu et al. 2022

131

### Open-Set Models

GroupViT: Semantic Segmentation Emerges from Text Supervision: Xu et al. 2022

132

### Open-Set Models

GroupVIT: Semantic Segmentation Emerges from Text Supervision: Xu et al. 2022

133

### Open-Set Models

GroupVIT: Semantic Segmentation Emerges from Text Supervision: Xu et al. 2022

134

### Bias in Vision and Language Models

Wrong	Right for the Right Reasons	Right for the Wrong Reasons	Right for the Right Reasons
Baseline: A <b>man</b> sitting at a desk with a laptop computer.	Our Model: A <b>woman</b> sitting in front of a laptop computer.	Baseline: A <b>man</b> holding a tennis racquet on a tennis court.	Our Model: A <b>man</b> holding a tennis racquet on a tennis court.

Women also Snowboard: Overcoming Bias in Caption Models: Burns et al. 2019

135

### Bias in Vision and Language Models

Women also Snowboard: Overcoming Bias in Caption Models: Burns et al. 2019

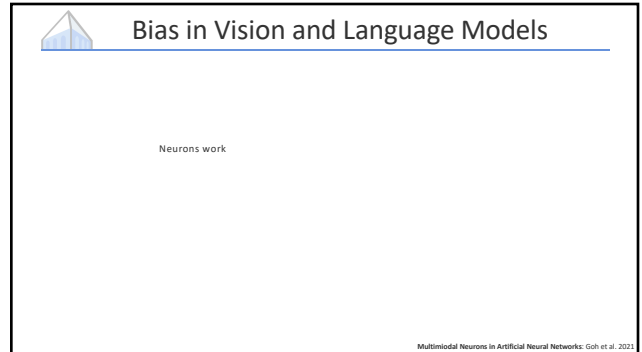
136

### 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

Evaluating CLIP: Towards Characterization of Broader Capabilities and Downstream Implications: Agarwal et al. 2021

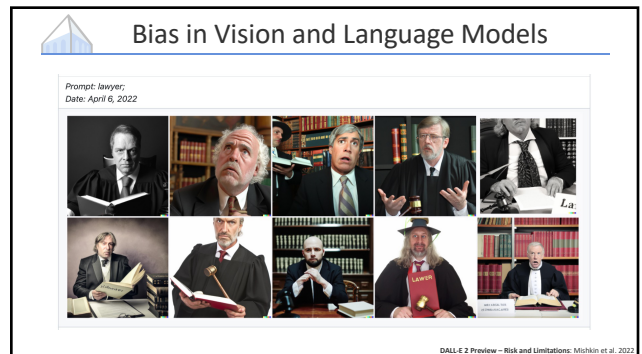
137



138



139



140



141



142