# Deep Learning for Images

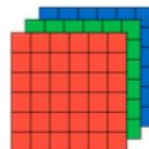**The Problem:** Semantic Gap



```
[[105 112 108 111 104  99 106  99  96 103 112 119 104  97  93  87]
 [ 91  98 102 106 104  79  98 103  99 105 123 136 110 105  94  85]
 [ 76  85  90 105 128 105  87  96  95  99 115 112 106 103  99  85]
 [ 99  81  81  93 120 131 127 100  95  98 102  99  96  93 101  94]
 [106  91  61  64  69  91  88  85 101 107 109  98  75  84  96  95]
 [114 108  85  55  55  69  64  54  64  87 112 129  98  74  84  91]
 [133 137 147 103  65  81  80  65  52  54  74  84 102  93  85  82]
 [128 137 144 140 109  95  86  70  62  65  63  63  60  73  86 101]
 [125 133 148 137 119 121 117  94  65  79  80  65  54  64  72  98]
 [127 125 131 147 133 127 126 131 111  96  89  75  61  64  72  84]
 [115 114 109 123 150 148 131 118 113 109 100  92  74  65  72  78]
 [ 89  93  90  97 108 147 131 118 113 114 113 109 106  95  77  80]
 [ 63  77  86  81  77  79 102 123 117 115 117 125 125 130 115  87]
 [ 62  65  82  89  78  71  80 101 124 126 119 101 107 114 131 119]
 [ 63  65  75  88  89  71  62  81 120 138 135 105  81  98 110 118]
 [ 87  65  71  87 106  95  69  45  76 130 126 107  92  94 105 112]
 [118  97  82  86 117 123 116  66  41  51  95  93  89  95 102 107]
 [164 146 112  80  82 120 124 104  76  48  45  66  88 101 102 109]
 [157 170 157 120  93  86 114 132 112  97  69  55  70  82  99  94]
 [130 128 134 161 139 100 109 118 121 134 114  87  65  53  69  86]
 [128 112  96 117 150 144 120 115 104 107 102  93  87  81  72  79]
 [123 107  96  86  83 112 153 149 122 109 104  75  80 107 112  99]
 [122 121 102  80  82  86  94 117 145 148 153 102  58  78  92 107]
 [122 164 148 103  71  56  78  83  93 103 119 139 102  61  69  84]]
```
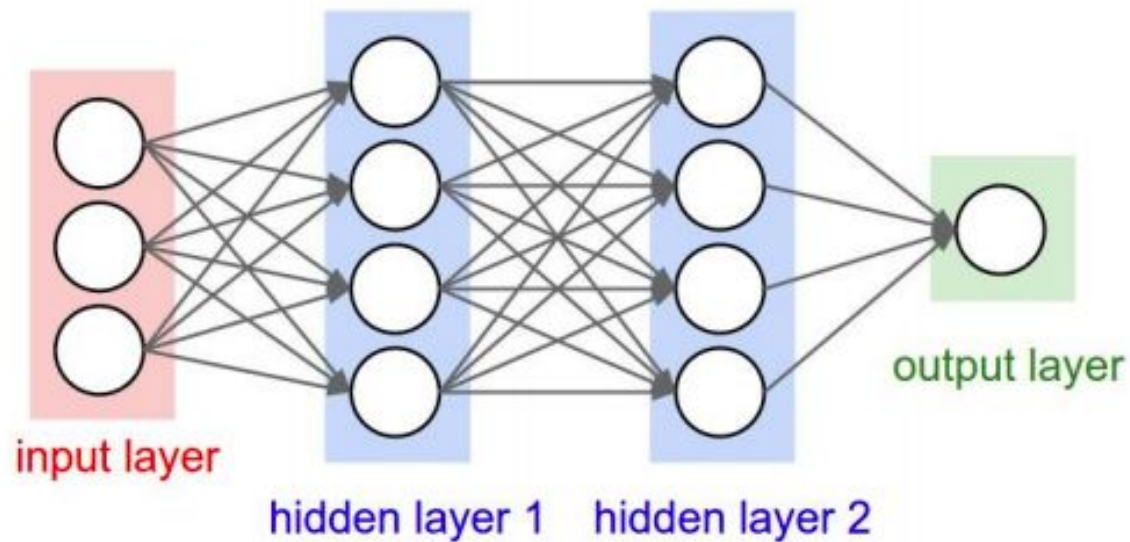
What the computer sees

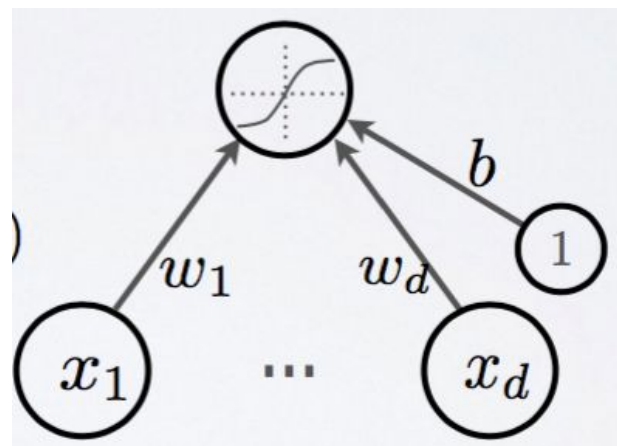An image is just a big grid of numbers between [0, 255]:

Example:
800 x 600 x 3
(3 channels RGB)

h
x
w
x
3

input layer

hidden layer 1    hidden layer 2
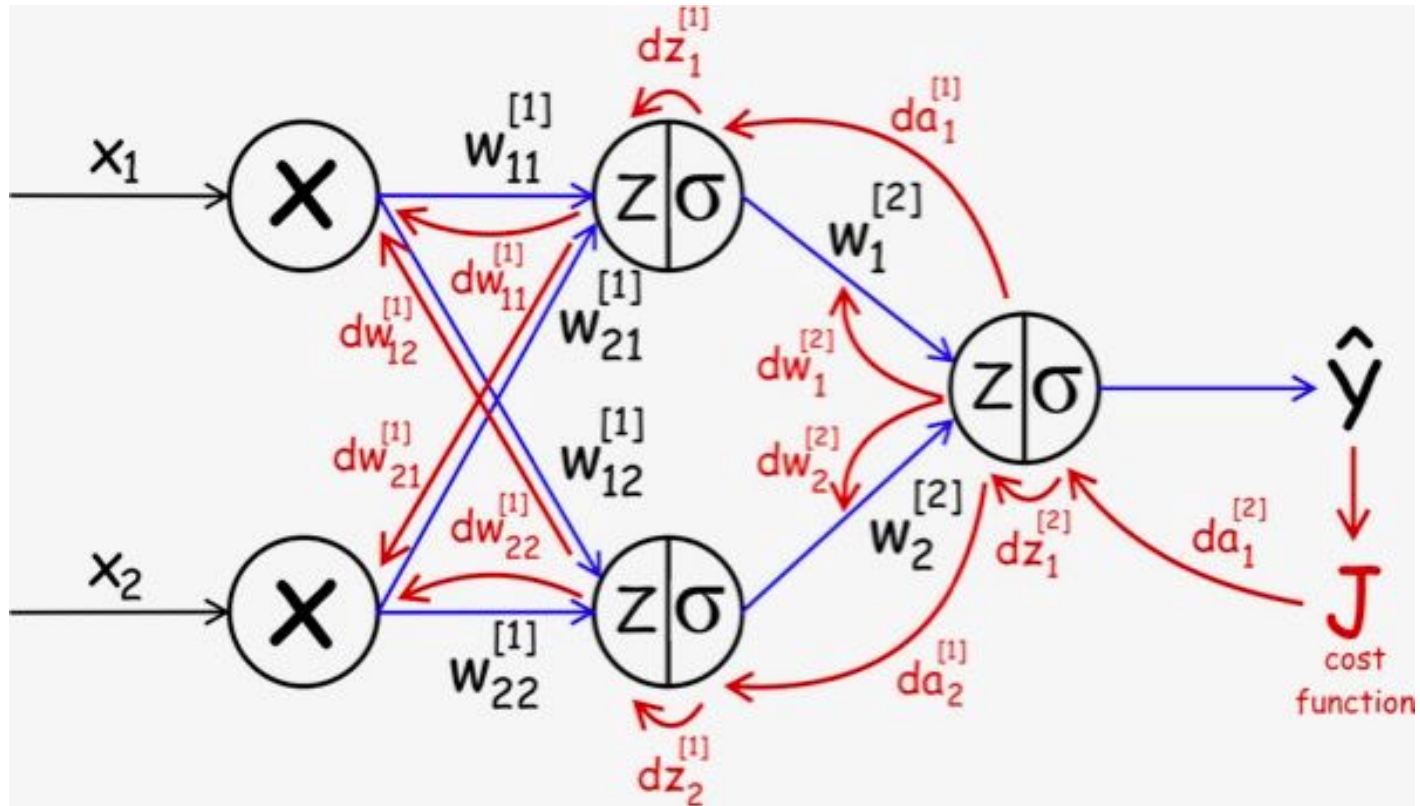
output layer

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

$$a(\mathbf{x}) = b + \sum_i w_i x_i$$

Suppose: 3 training examples, 3 classes.
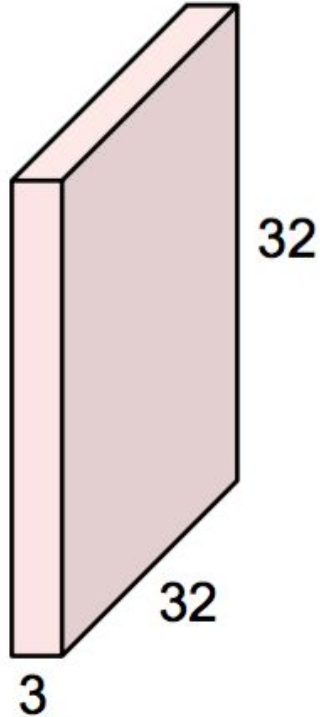With some W the scores $f(x, W) = Wx$ are:



| | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

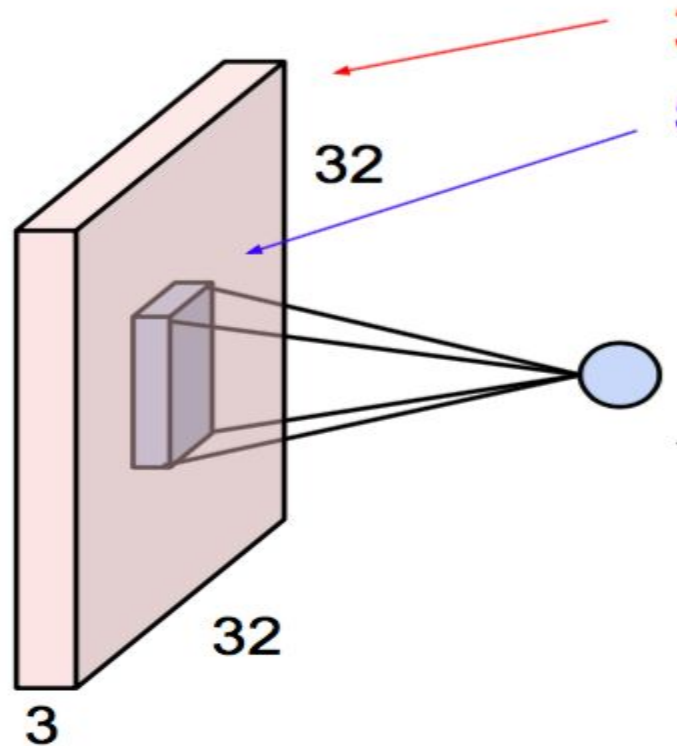# Backpropagation Algorithm
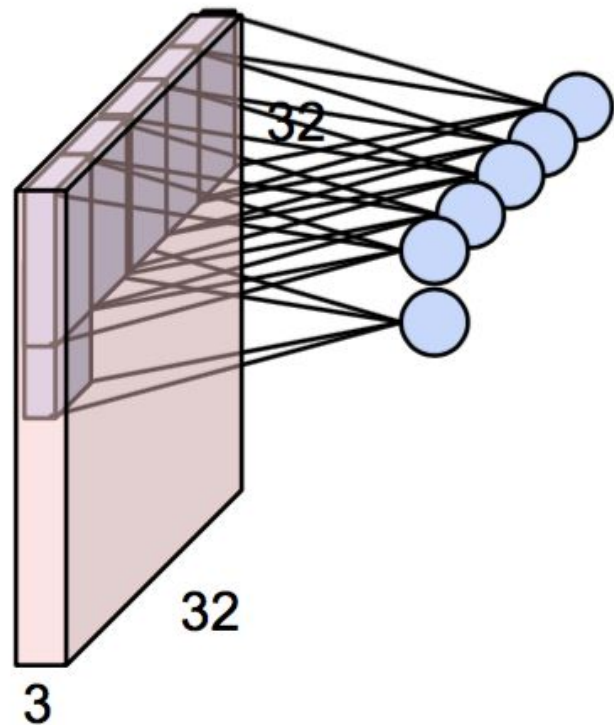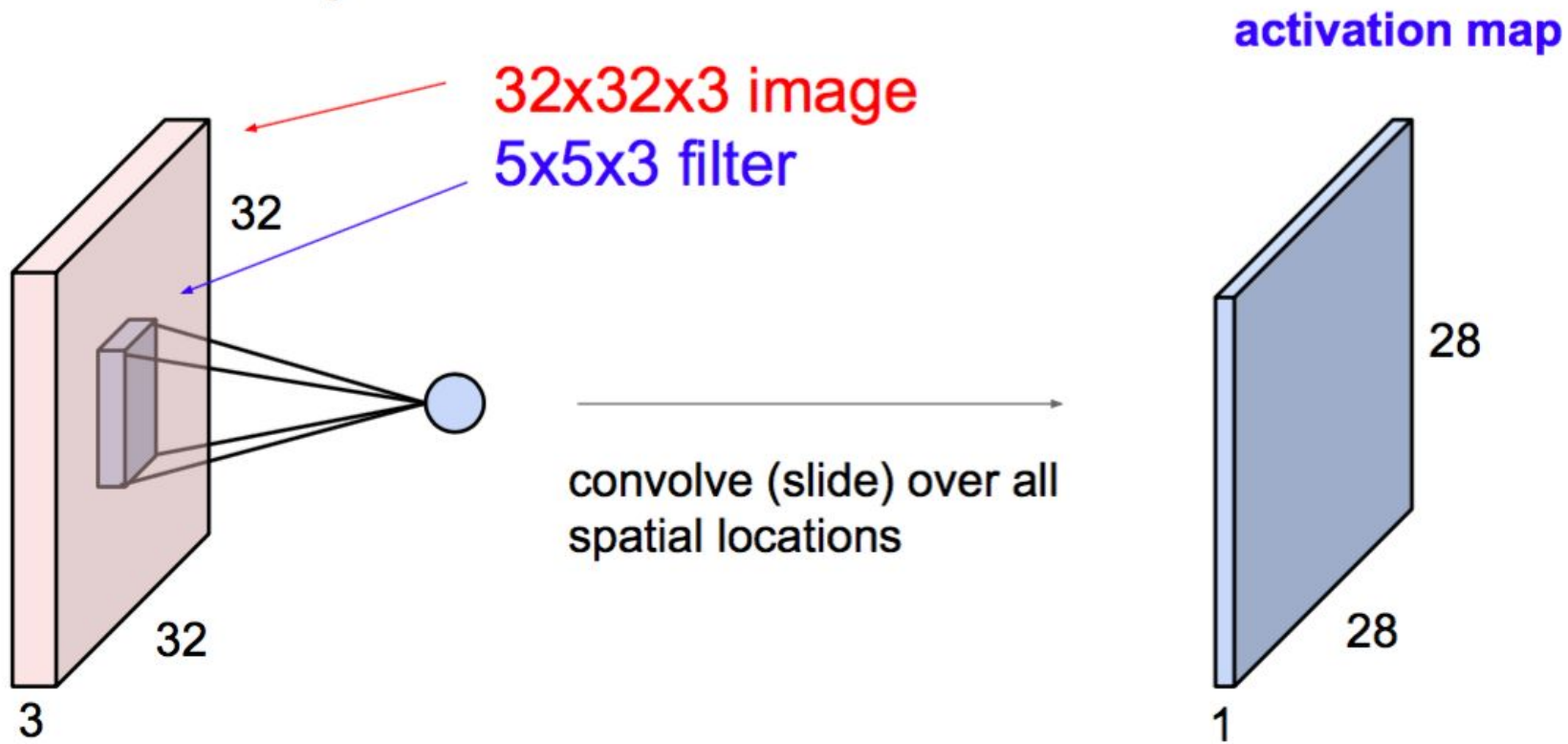
# Convolution

32x32x3 image
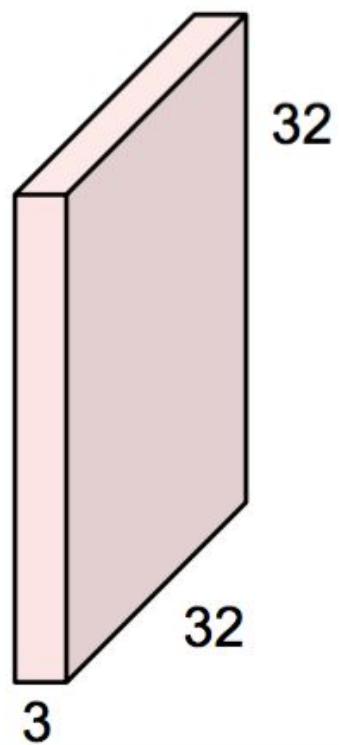


32

32

3

5x5x3 filter



**Co**

i.e.

co

## Convolution Layer

$$w^T x + b$$
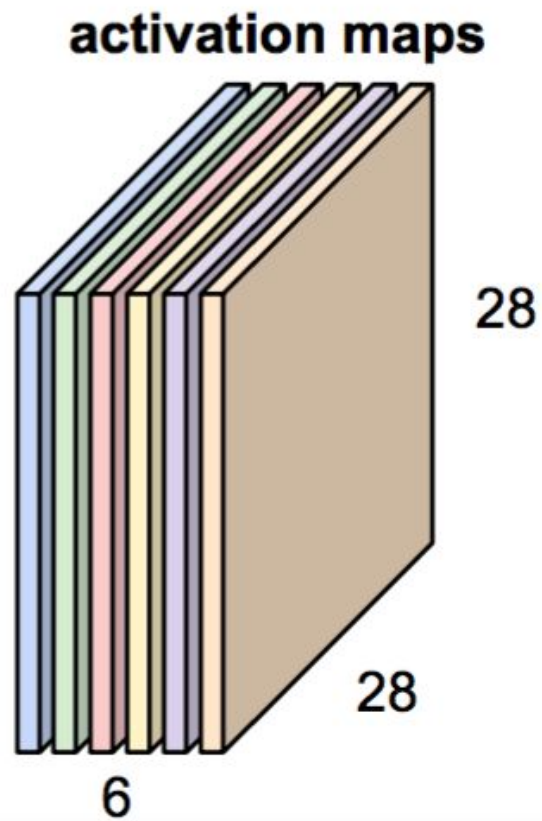
32

32

3

32

32

3

**32x32x3 image**

**5x5x3 filter**

32

32

3

convolve (slide) over all
spatial locations

**activation map**

28

28

1

**activation maps**

32

32

3

Convolution Layer
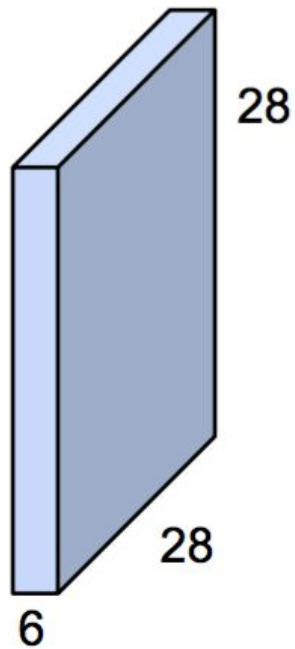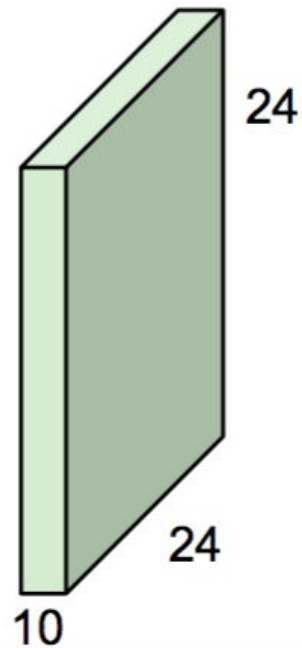
28

28

6

32
32
3

CONV,
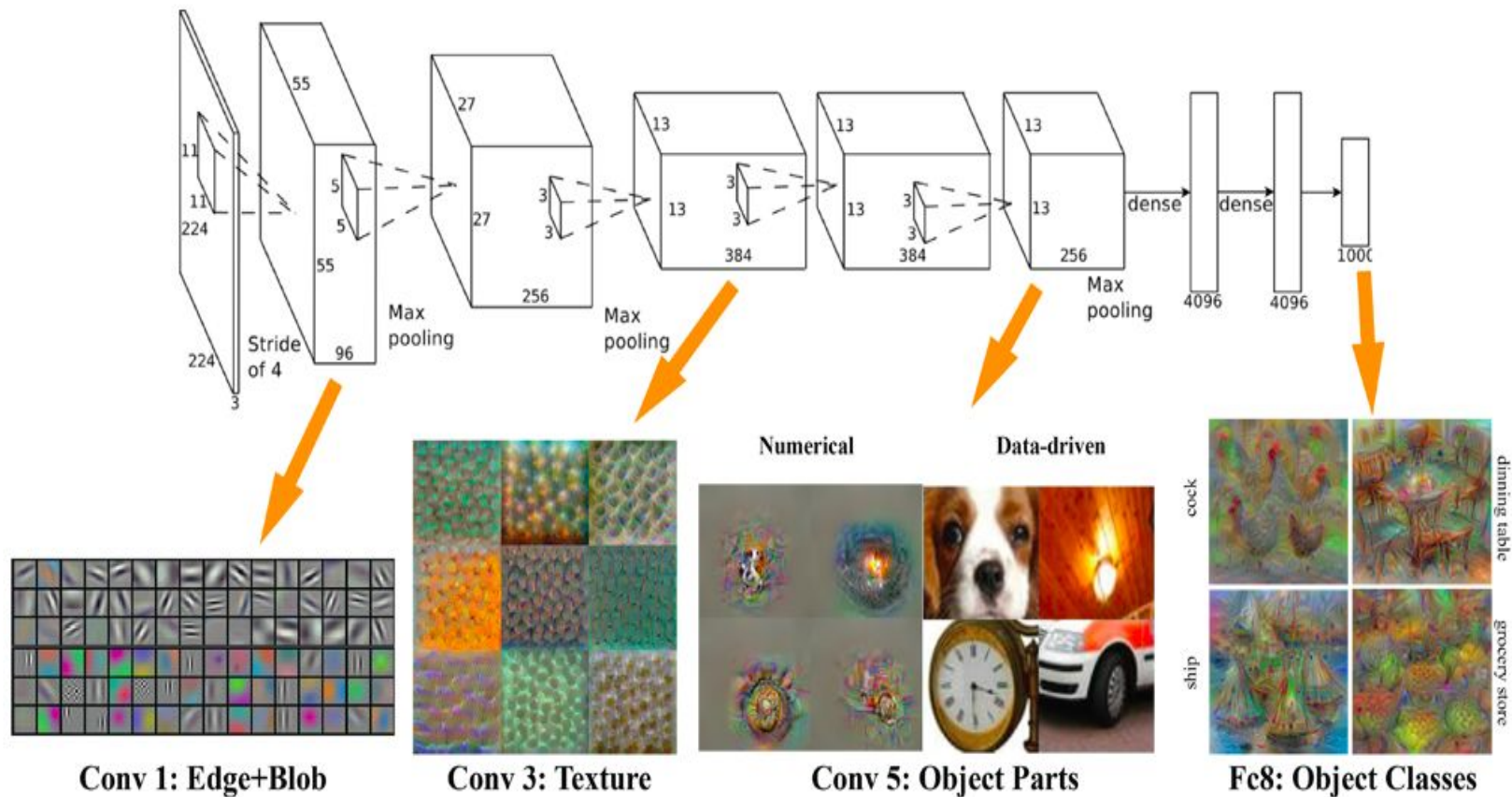ReLU
e.g. 6
5x5x3
filters

28
28
6

CONV,
ReLU
e.g. 10
5x5x**6**
filters

24
24
10

CONV,
ReLU

....

AlexNet / VGG-F network visualized by **mNeuron**.

# Preview

*[Zeiler and Fergus 2013]*



Low-level features → Mid-level features → High-level features → Linearly separable classifier

VGG-16 Conv1_1

VGG-16 Conv3_2

VGG-16 Conv5_3

# Deep Learning Consists of two characteristics

1. End to End Training due to Backpropagation Algorithm.
2. Hierarchical Learning of Features.

Deep Learning become popular recently due to three reasons :-

1. More Data
2. More Computational Power due to GPU(Parallel Processing)
3. Better Algorithms

mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

Imagenet Challenge(Dataset)- 1.2 million images and 1000 categories

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

# Transfer learning: idea



Source labels

Source model

Large amount of data/labels

Source data
*E.g. ImageNet*

Transfer Learned Knowledge

Target labels

Target model

Small amount of data/labels

Target data
*E.g. PASCAL*

# At an abstract level, deep neural networks operate with some similar principals to the real brain (though there are some important differences!)



Spatial convolution over image input

Yamins & DiCarlo (2016), *Nature Neuroscience*

# Computer Vision Tasks

| **Classification** | **Semantic Segmentation** | **Object Detection** | **Instance Segmentation** |
|---|---|---|---|



**CAT**

**GRASS, CAT, TREE, SKY**

**DOG, DOG, CAT**

**DOG, DOG, CAT**

No spatial extent

No objects, just pixels

Multiple Object

# Image Segmentation



Input:
3 x H x W

Conv → Conv → Conv → Conv → argmax

Convolutions:
D x H x W

Scores:
C x H x W

Predictions:
H x W

**Downsampling:** Pooling, strided convolution

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

**Upsampling:** ???

Input: $3 \times H \times W$

High-res: $D_1 \times H/2 \times W/2$

Med-res: $D_2 \times H/4 \times W/4$

Low-res: $D_3 \times H/4 \times W/4$

Med-res: $D_2 \times H/4 \times W/4$

High-res: $D_1 \times H/2 \times W/2$

Predictions: $H \times W$

# Object Detection

# R-CNN



SVMs — Classify regions with SVMs

Forward each region through ConvNet

Warped image regions (224x224 pixels)

Regions of Interest (RoI) from a proposal method (~2k)

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Objects + <u>Relationships</u> = Scene Graphs



108,077 Images
5.4 Million Region Descriptions
1.7 Million Visual Question Answers
3.8 Million Object Instances
2.8 Million Attributes
2.3 Million Relationships
Everything Mapped to Wordnet Synsets

Krishna, Ranjay, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations." International Journal of Computer Vision 123, no. 1 (2017): 32-73.

# Scene Graph Prediction

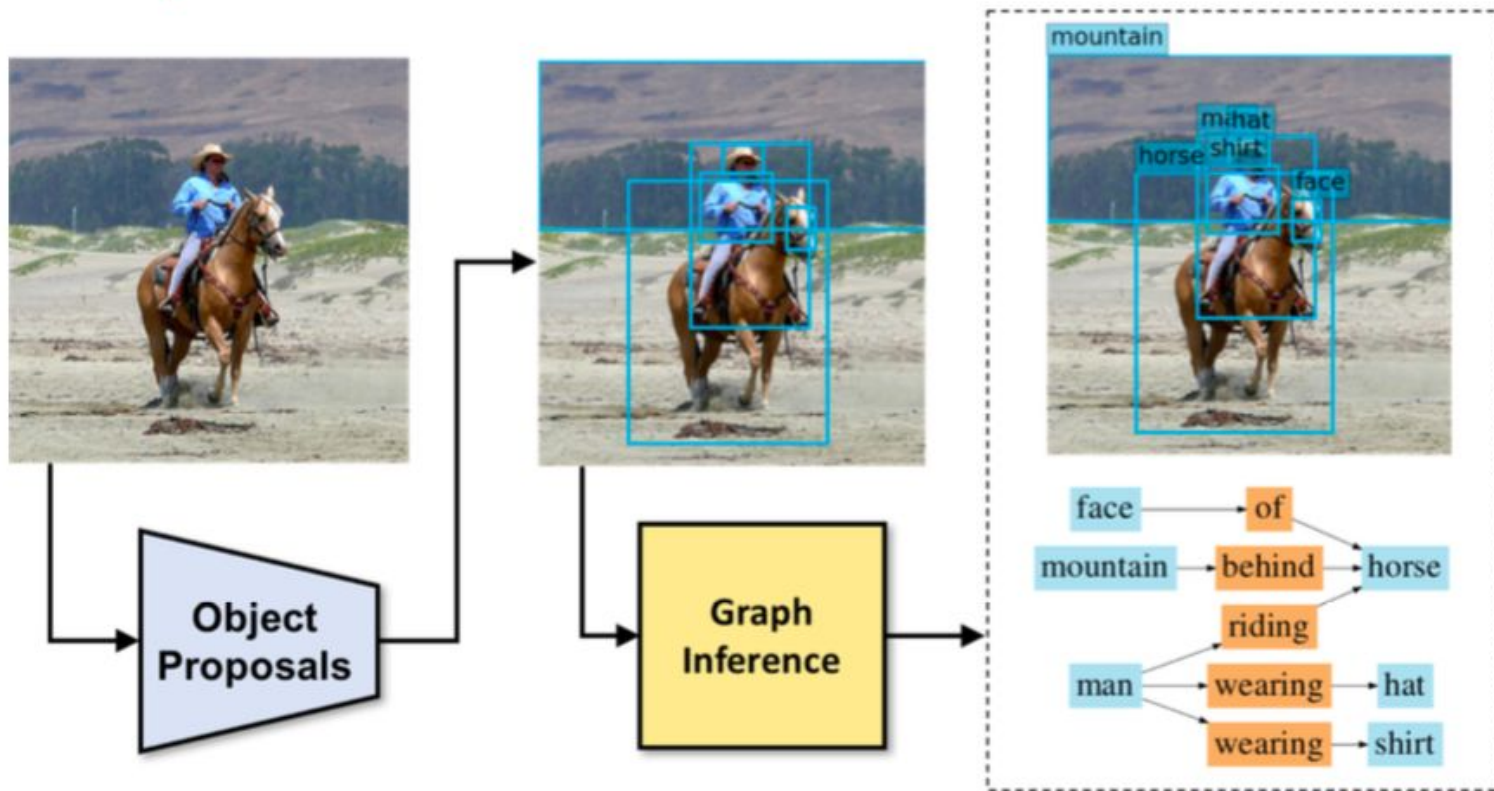# Language and vision

## Captioning



"man in black shirt is playing guitar."

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

## Dense Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

## Referring Expressions

*largest elephant standing behind baby elephant.*



Zhang, Liu, and Chang, "Grounding Referring Expressions in Images by Variational Context", CVPR 2018

# Visual Question Answering (VQA)

**VQA**



What color are her eyes?
What is the mustache made of?

How many slices of pizza are there?
Is this a vegetarian pizza?

Is this person expecting company?
What is just under the tree?

Does it appear to be rainy?
Does this person have 20/20 vision?

- Understanding of visual input
- Understanding of language
- World knowledge
- Reasoning

- Can ask about anything
- Easier to evaluate (at least for multiple choice)

"VQA: Visual Question Answering"
[Agrawal et al, ICCV 2015]

# Visual Question Answering (VQA)

## VQA

Who is wearing glasses?
man          woman
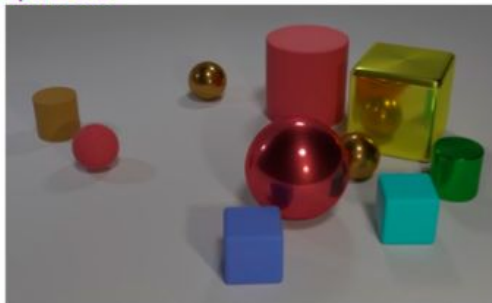


Is the umbrella upside down?
yes          no



"Making the V in VQA Matter: Elevating the Role of Image
Understanding in Visual Question Answering"
[Goyal et al, CVPR 2017]

## CLEVR

Questions in CLEVR test various aspects of visual
reasoning including attribute identification, counting,
comparison, spatial relationships, and logical
operations.



Q: Are there an equal number of large things and metal spheres?
Q: What size is the cylinder that is left of the brown metal thing
that is left of the big sphere?
Q: There is a sphere with the same size as the metal cube; is it
made of the same material as the small red sphere?
Q: How many objects are either small cylinders or red things?

"CLEVR: A Diagnostic Dataset for Compositional Language
and Elementary Visual Reasoning"
[Johnson et al, CVPR 2017]

## GQA



Is the bowl to the right of the green apple?
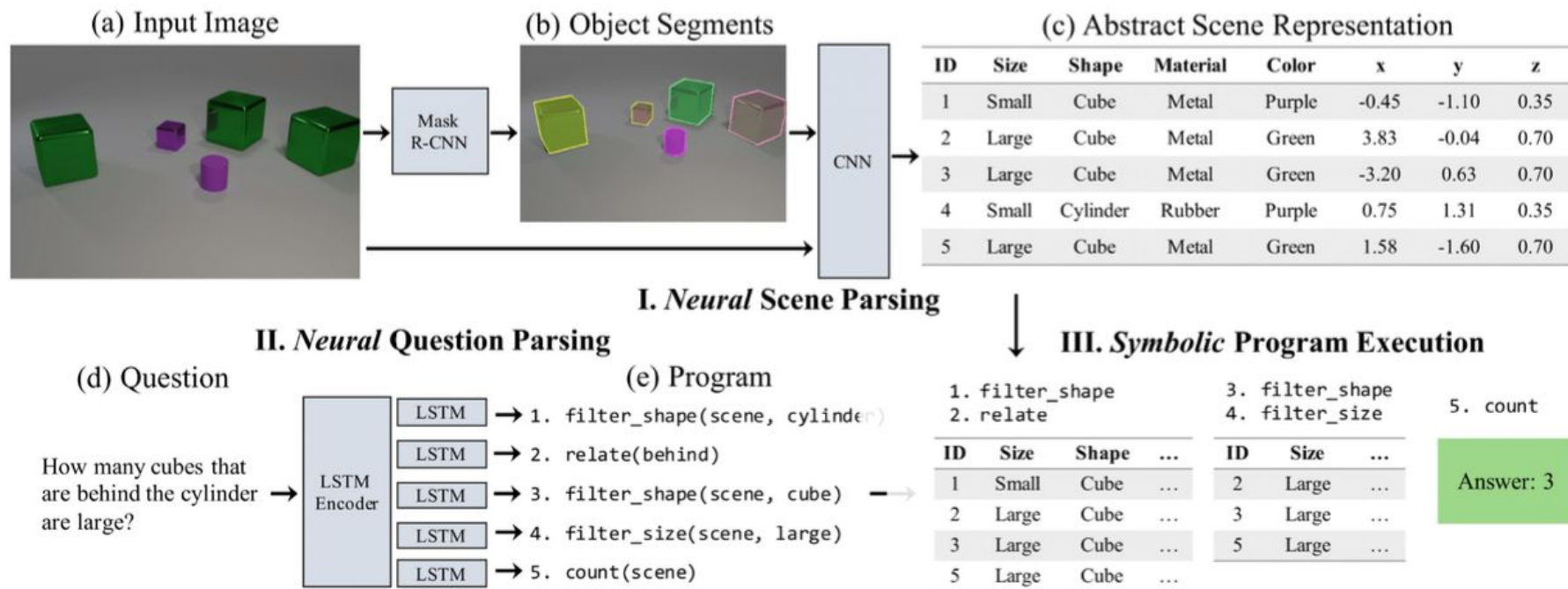What type of fruit in the image is round?
What color is the fruit on the right side, red or green?
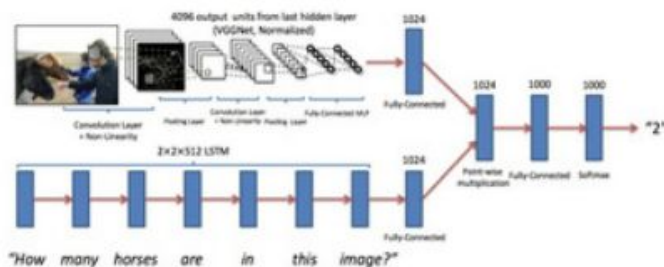Is there any milk in the bowl to the left of the apple?

"GQA: A New Dataset for Real-World Visual Reasoning and
Compositional Question Answering"
[Hudson and Manning, CVPR 2019]

# Reasoning



(a) Input Image

(b) Object Segments

Mask R-CNN

CNN

(c) Abstract Scene Representation

| ID | Size | Shape | Material | Color | x | y | z |
|---|---|---|---|---|---|---|---|
| 1 | Small | Cube | Metal | Purple | -0.45 | -1.10 | 0.35 |
| 2 | Large | Cube | Metal | Green | 3.83 | -0.04 | 0.70 |
| 3 | Large | Cube | Metal | Green | -3.20 | 0.63 | 0.70 |
| 4 | Small | Cylinder | Rubber | Purple | 0.75 | 1.31 | 0.35 |
| 5 | Large | Cube | Metal | Green | 1.58 | -1.60 | 0.70 |

**I. *Neural* Scene Parsing**

**II. *Neural* Question Parsing**

**III. *Symbolic* Program Execution**

(d) Question

How many cubes that are behind the cylinder are large?

LSTM Encoder

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

(e) Program

1. filter_shape
2. relate

| ID | Size | Shape | ... |
|---|---|---|---|
| 1 | Small | Cube | ... |
| 2 | Large | Cube | ... |
| 3 | Large | Cube | ... |
| 5 | Large | Cube | ... |

3. filter_shape
4. filter_size

| ID | Size | ... |
|---|---|---|
| 2 | Large | ... |
| 3 | Large | ... |
| 5 | Large | ... |

5. count

Answer: 3

"Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding"
[Yi, Wu, Gan, Torralba, Kohli, and Tennebaum, NeurIPS 2018]

# Task- and dataset-specific models
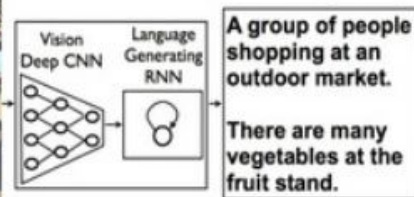


Visual Question Answering [Antol et. al. 2015]
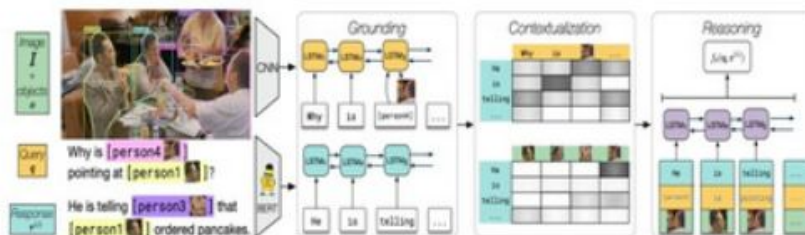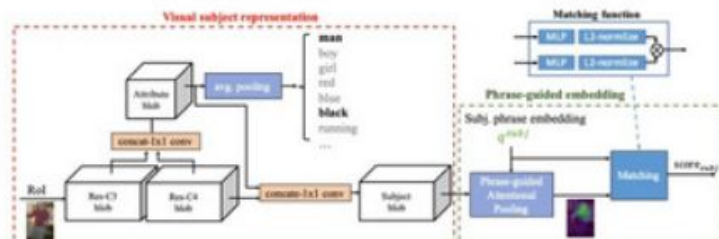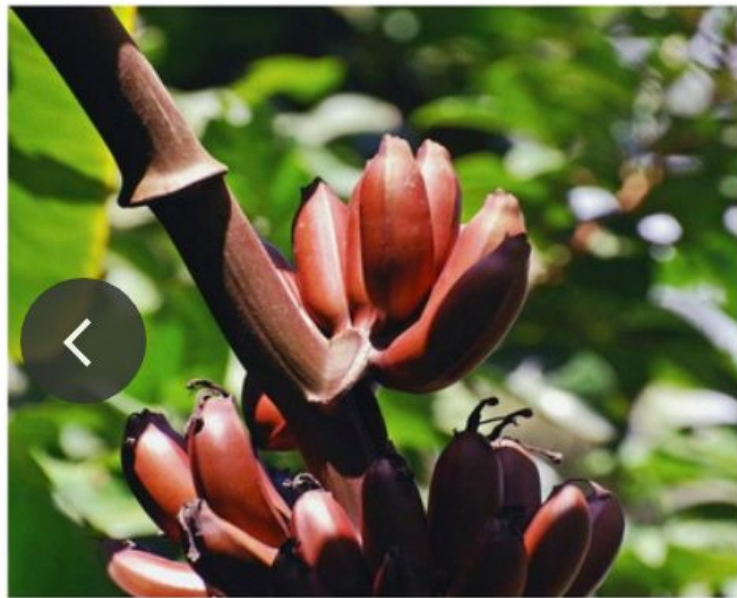


Image Captioning [Vinyals et. al. 2015]



Visual Commonsense Reasoning [Zellers et. al. 2018]



Refer Expression [Yu et. al 2018]

# Task- and dataset-specific models



VQA model:
Q: What type of plant is this?
A: Banana

Captioning model:
A bunch of red and yellow flowers on a branch.

Common model for visual grounding
Leverage for a variety of vision-and-language tasks

# ViLBERT Multi-Task [CVPR 2020]

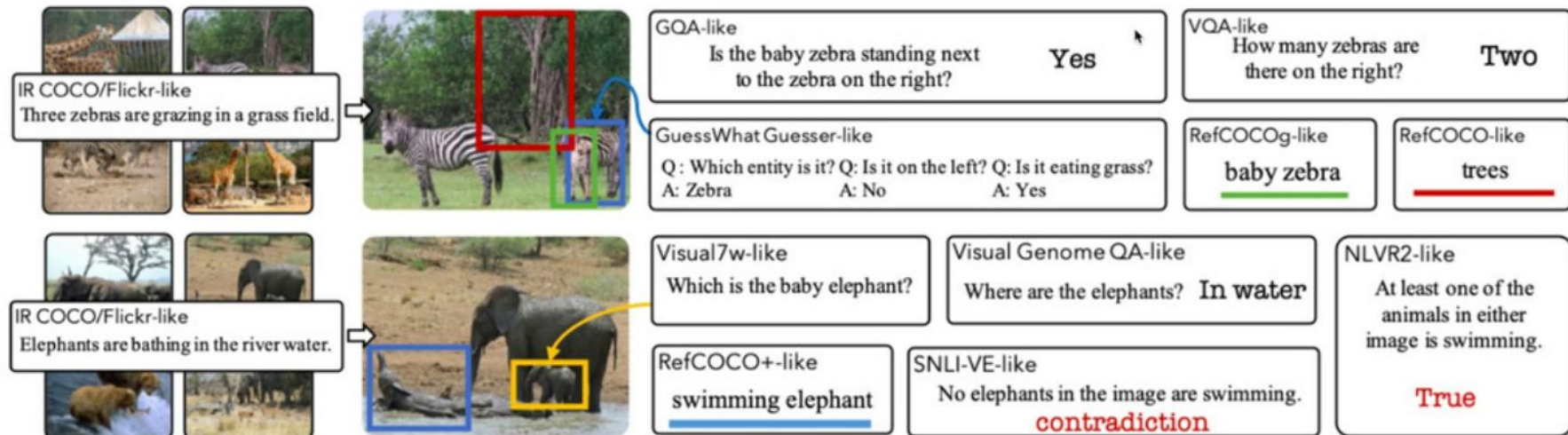1 model for 12 tasks!

Higher performance, 1/12th the model size!

SOTA on 7 after fine-tuning

Jiasen Lu

Vedanuj Goswami

IR COCO/Flickr-like
Three zebras are grazing in a grass field.

GQA-like
Is the baby zebra standing next to the zebra on the right?   Yes

VQA-like
How many zebras are there on the right?   Two

GuessWhat Guesser-like
Q : Which entity is it? Q: Is it on the left? Q: Is it eating grass?
A: Zebra          A: No          A: Yes

RefCOCOg-like
baby zebra

RefCOCO-like
trees

IR COCO/Flickr-like
Elephants are bathing in the river water.

Visual7w-like
Which is the baby elephant?

Visual Genome QA-like
Where are the elephants?   In water

NLVR2-like
At least one of the animals in either image is swimming.

RefCOCO+-like
swimming elephant

SNLI-VE-like
No elephants in the image are swimming.
contradiction

True

92

# Extending to 3D

| Classification | Semantic Segmentation | Object Detection | Instance Segmentation |
|---|---|---|---|

**Apartment**

**Table**, **Bed**, **Couch**, **Cabinet**

**Bed**, **Couch**, **Cabinet**, **Cabinet**

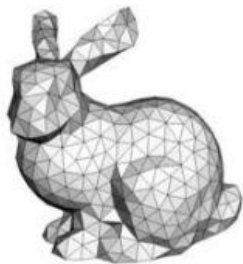**Bed**, **Couch**, **Cabinet**, **Desk**
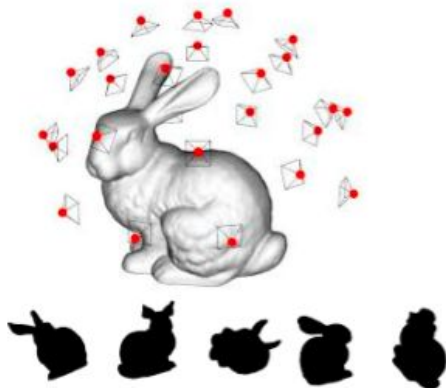
Single label

No objects

Multiple Object

"ScanNet: Richly-annotated 3D Reconstructions of Indoor Scenes"
[Dai et al, CVPR 2017]

82

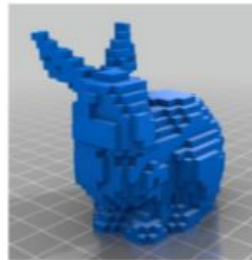# Extending to 3D - Representation
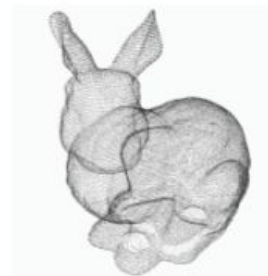
**Surface:**
Triangle Mesh

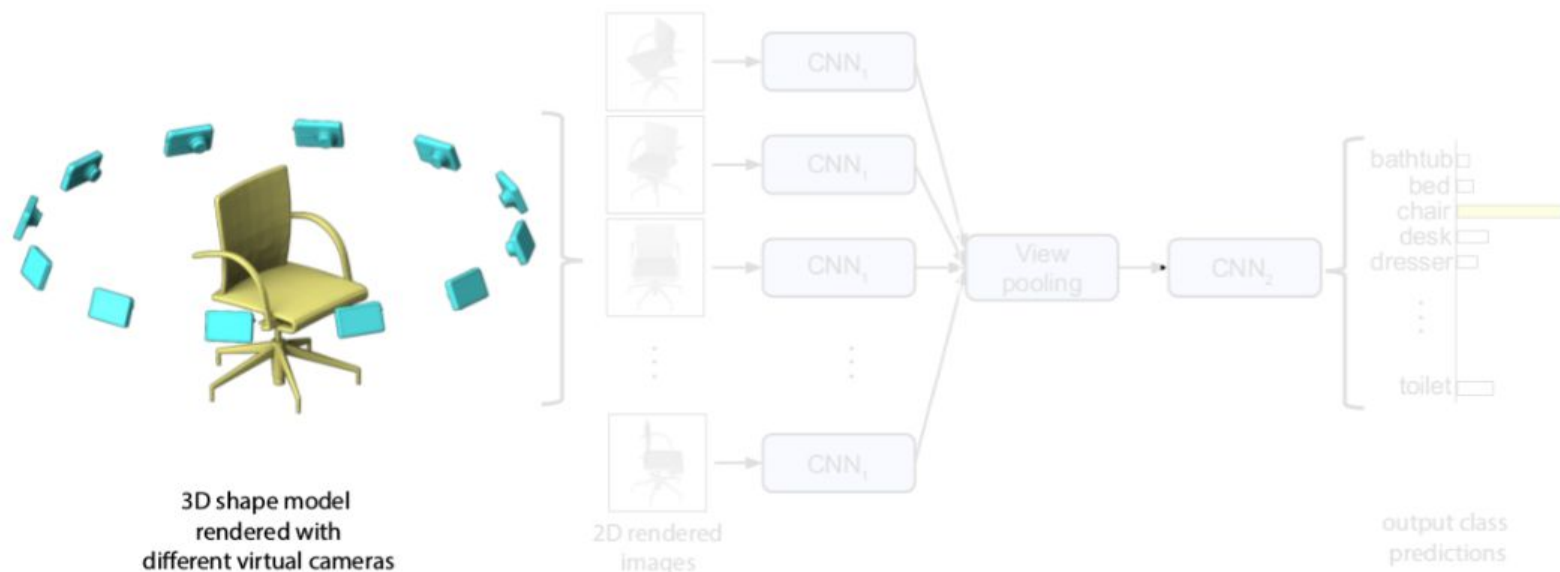**Multi-View:**
Set of Images

**Volumetric:**
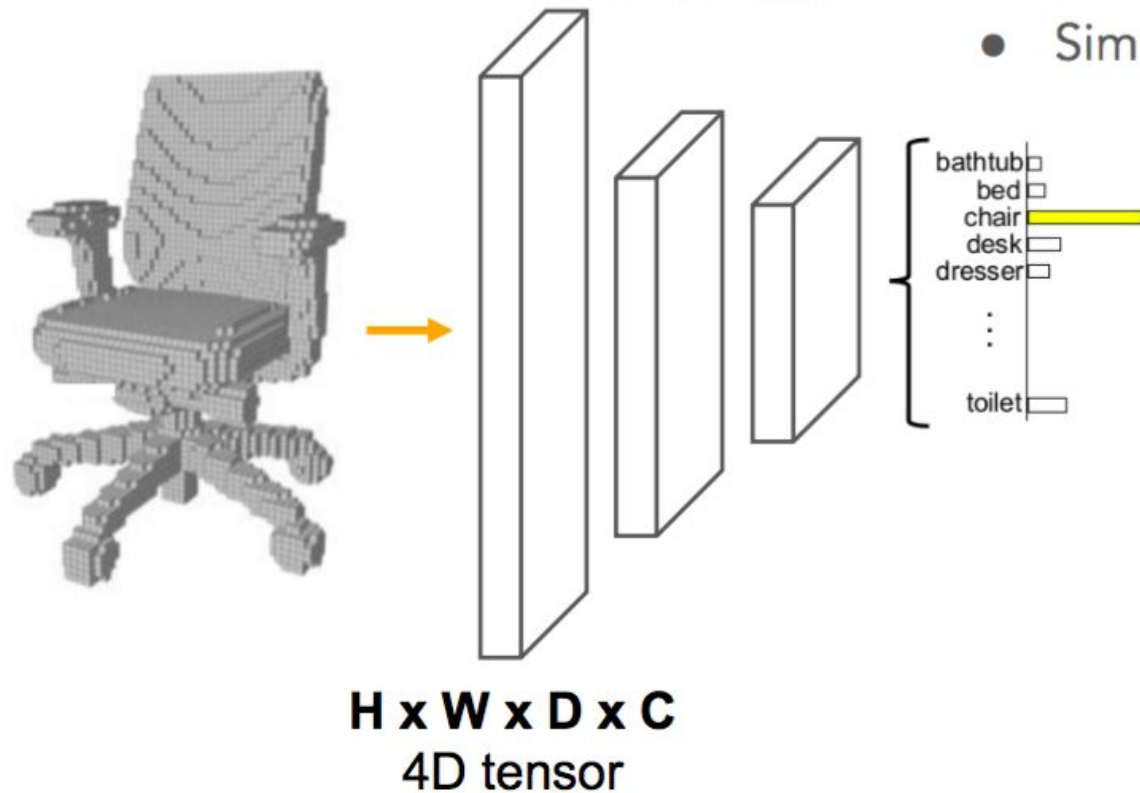Voxels

**Pointcloud:**
Set of points

# Multiview



"Multi-view Convolutional Neural Networks for 3D Shape Recognition"
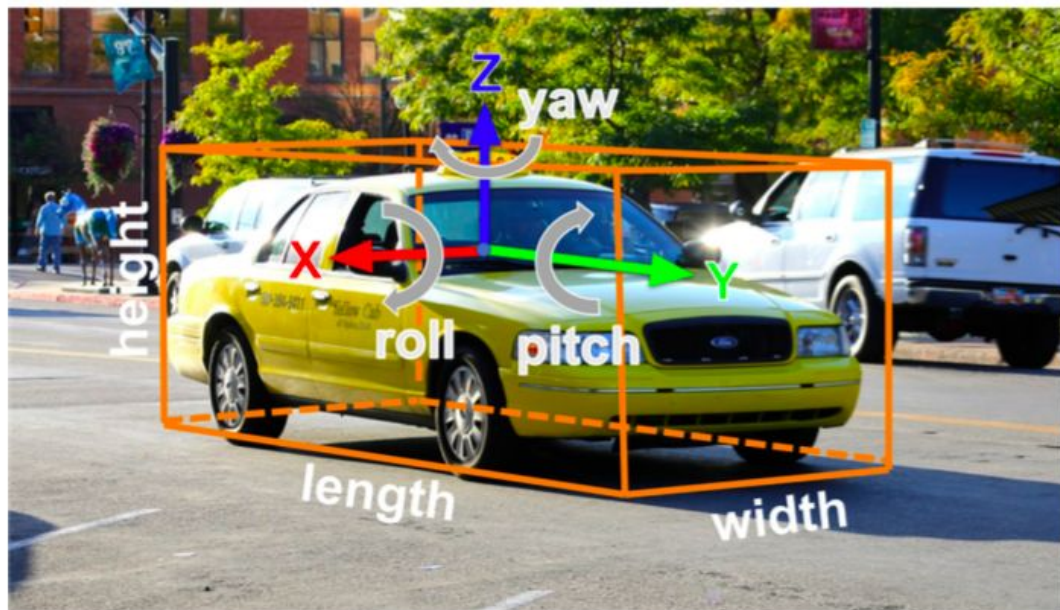[Su, Maji, Kalogerakis, Learned-Miller, ICCV 2015]

# Voxels

Convolve as before!

- Simple extension of pixel

bathtub
bed
chair
desk
dresser

⋮

toilet

**H x W x D x C**
4D tensor

# 3D Object Detection



2D Object Detection:
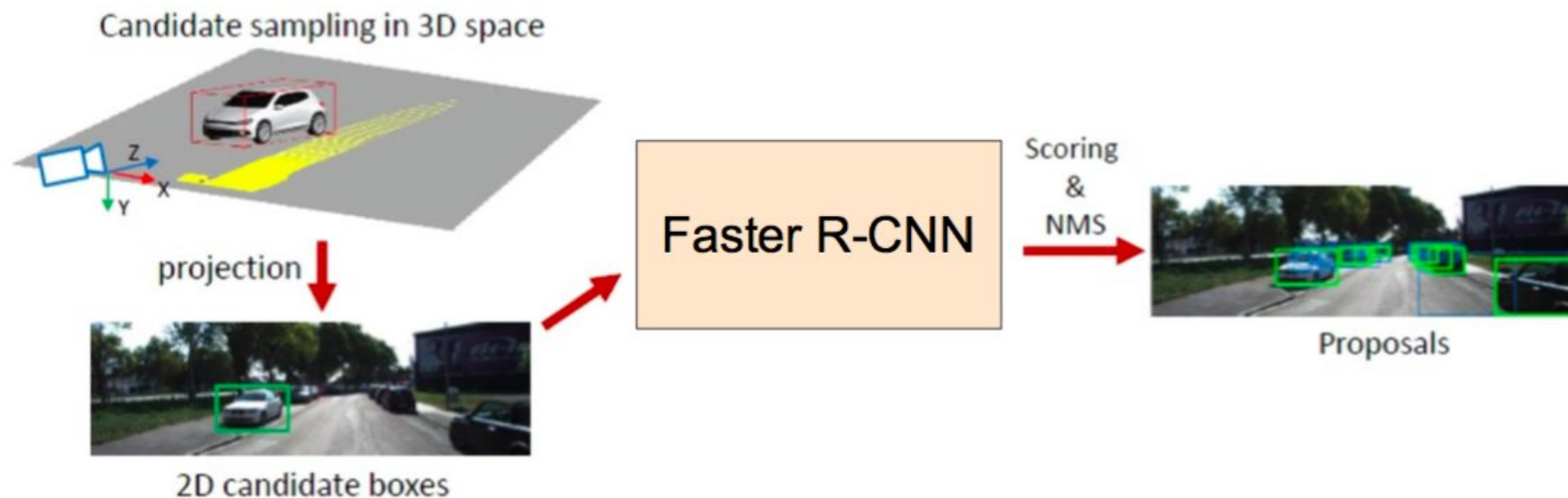2D bounding box
 (x, y, w, h)

3D Object Detection:
3D oriented bounding box
 (x, y, z, w, h, l, r, p, y)

Simplified bbox: no roll & pitch
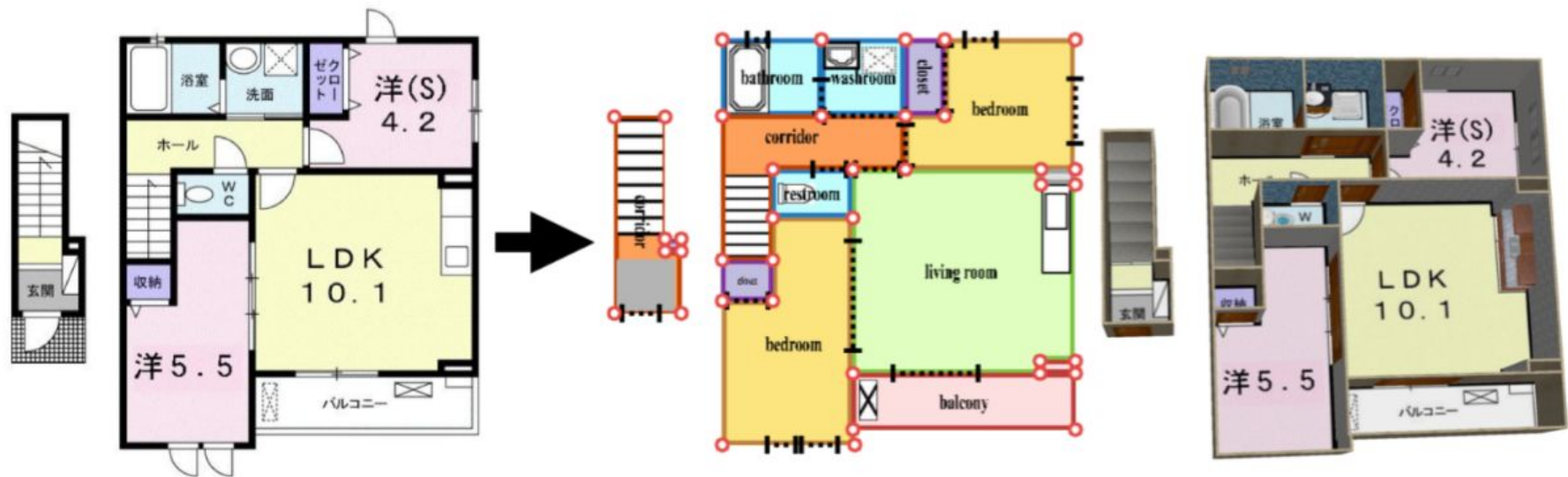
Much harder problem than 2D object detection!
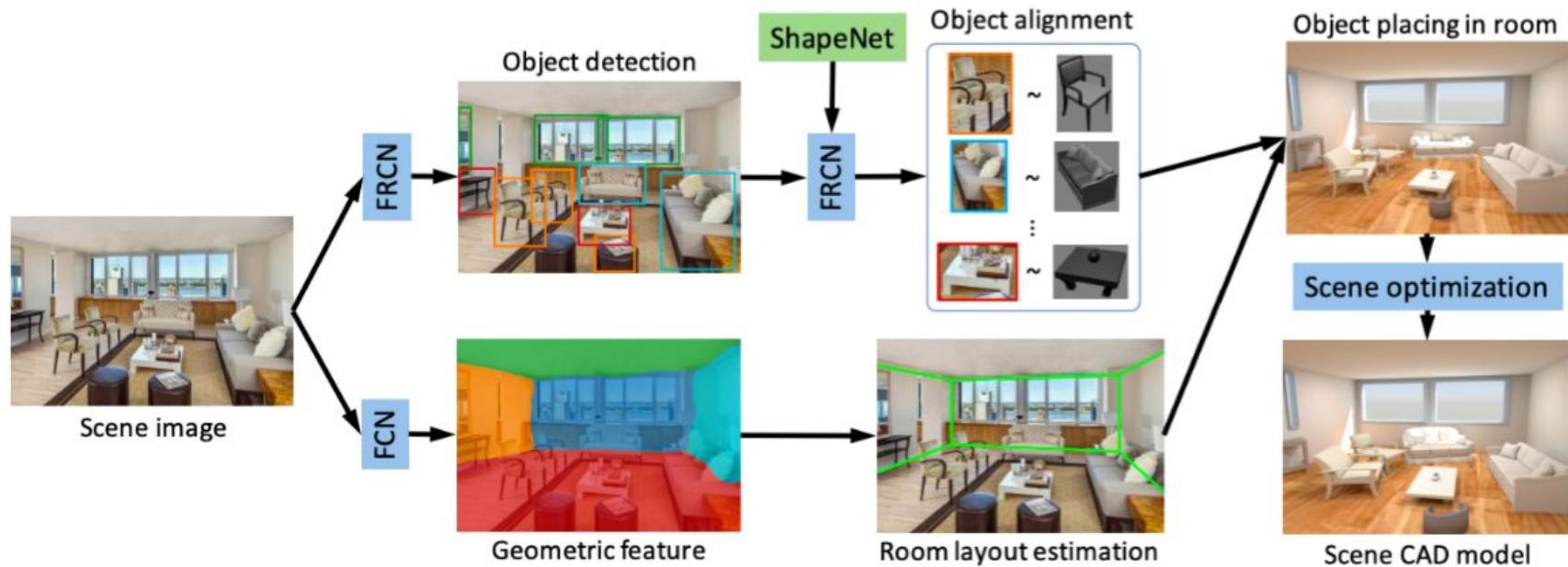
# 3D Object Detection: Monocular Camera



Candidate sampling in 3D space

projection

2D candidate boxes

Faster R-CNN

Scoring & NMS

Proposals

- Same idea as Faster RCNN, but proposals are in 3D
- 3D bounding box proposal, regress 3D box parameters + class score

Chen, Xiaozhi, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. "Monocular 3d object detection for autonomous driving." CVPR 2016.

# 2D floorplan to 3D model



"Raster-to-Vector: Revisiting Floorplan Transformation"
[Liu, Wu, Kohli, Furukawa, ICCV 2017]

"IM2CAD"
[Izadinia et al, CVPR 2017]

114

# Can we generate 3D scenes from (almost) scratch?



Empty Room

Select and arrange
3D models

→ Nicely arranged
Living Room

3D model database