# **CS 4644-DL / 7643-A: LECTURE 21 DANFEI XU**

- 2D Computer Vision (Continued)
- 3D Vision: Representations and Neural Rendering

## Computer Vision Tasks

#### Classification



**CAT** 

No spatial extent

Semantic Segmentation



TREE, SKY

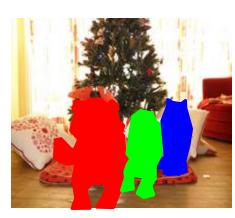
No objects, just pixels

Object Detection



DOG, DOG, CAT

## Instance Segmentation



DOG, DOG, CAT

Multiple Object

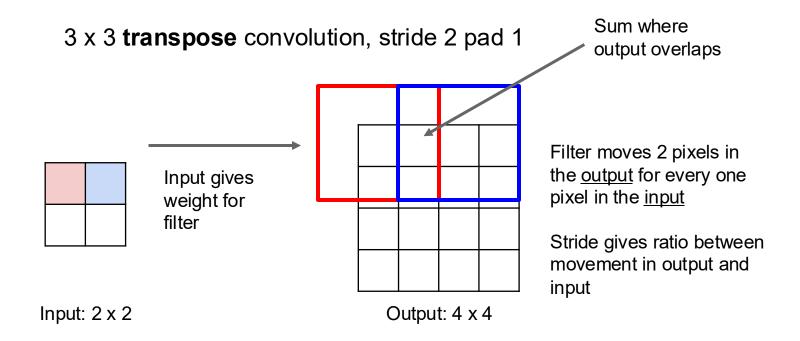
This image is CC0 public domain

## Semantic Segmentation Idea: Fully Convolutional

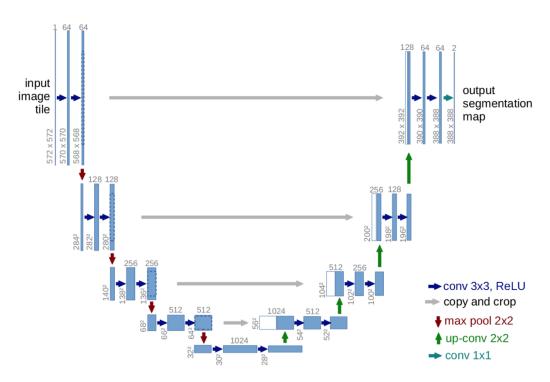
Design network as a bunch of convolutional layers, with **Downsampling: Upsampling:** downsampling and upsampling inside the network! Pooling, strided ??? convolution Med-res: Med-res:  $D_2 \times H/4 \times W/4$  $D_2 \times H/4 \times W/4$ Low-res:  $D_3 \times H/4 \times W/4$ Input: High-res: CxHxWHigh-res: Predictions: 3xHxWD<sub>1</sub> x H/2 x W/2  $D_1 \times H/2 \times W/2$ HxW

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

#### Learnable Upsampling: Transposed Convolution



## Semantic Segmentation: U-Net



Idea: Concatenate feature maps from the downsampling stage with the features in the upsampling stage.

Very commonly used today!

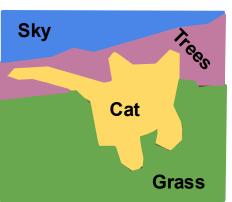
Ronneberger O, Fischer P, Brox T, 2015

### Semantic Segmentation

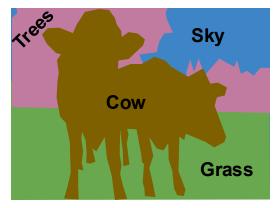
Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels









## **Object Detection**

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Instance Segmentation



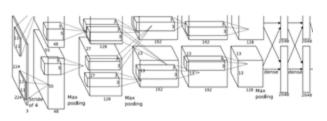
DOG, DOG, CAT

Multiple Object

## Object Detection: Multiple Objects

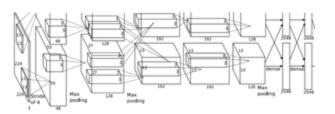
## Each image needs a different number of outputs!





CAT: (x, y, w, h) 4 numbers





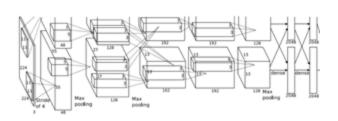
DOG: (x, y, w, h)

DOG: (x, y, w, h)

12 numbers

CAT: (x, y, w, h)





DUCK: (x, y, w, h) Many

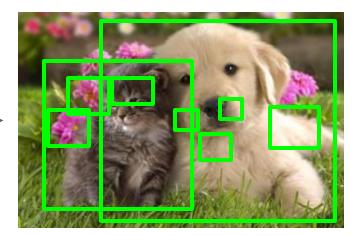
DUCK: (x, y, w, h) numbers!

. . . .

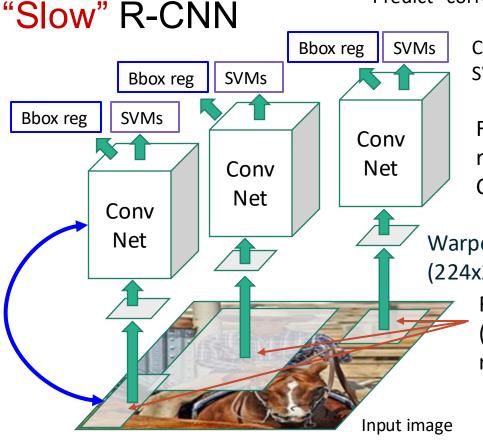
## Region Proposals: Selective Search

- Find "blobby" image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU





Predict "corrections" to the RoI: 4 numbers: (dx, dy, dw, dh)



Classify regions with SVMs

Forward each region through ConvNet

Warped image regions (224x224 pixels)

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

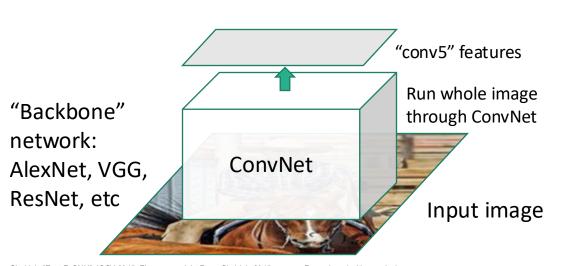
**Problem**: Very slow! Need to do ~2k independent forward passes for each image!

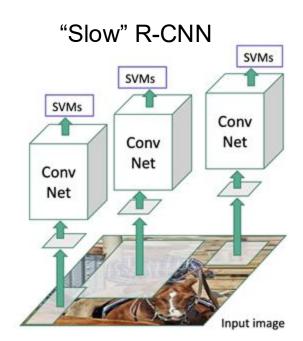
Idea: Pass the image through convnet before cropping! Crop the conv feature instead!

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

Figure copyright Ross Girshick, 2015;  $\underline{\text{source}}$ . Reproduced with permission.

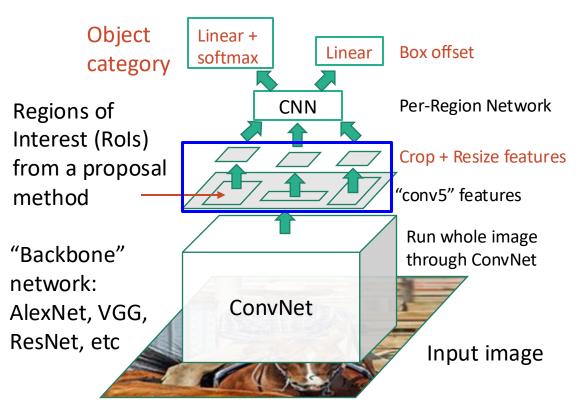
#### Fast R-CNN

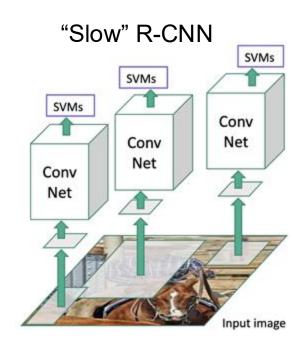




Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

#### Fast R-CNN





Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission

### Cropping Features: Rol Pool

grid cells Project proposal onto features **CNN** 

Input Image (e.g. 3 x 640 x 480)

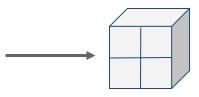
Image features: C x H x W (e.g. 512 x 20 x 15)

"Snap" to

Problem: Region features slightly misaligned

Divide into 2x2 grid of (roughly) equal subregions

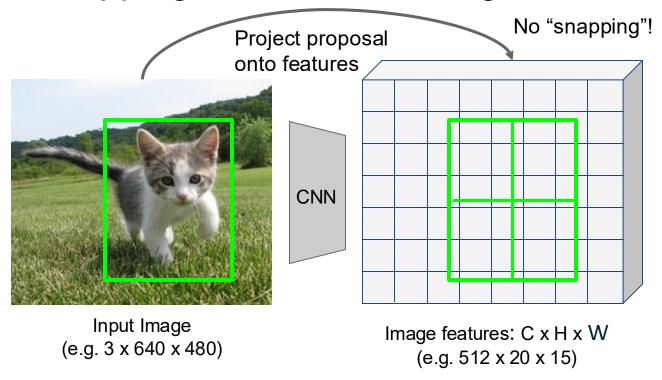
Max-pool within each subregion



Region features (here 512 x 2 x 2; In practice e.g 512 x 7 x 7)

Region features always the same size even if input regions have different sizes!

Girshick, "Fast R-CNN", ICCV 2015.



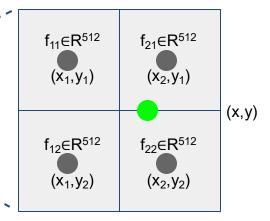
No "snapping"! Project proposal onto features **CNN** Input Image Image features: C x H x W (e.g. 3 x 640 x 480) (e.g. 512 x 20 x 15)

Sample at regular points in each subregion using bilinear interpolation

He et al, "Mask R-CNN", ICCV 2017

No "snapping"! Project proposal onto features **CNN** Input Image Image features: C x H x W

Sample at regular points in each subregion using bilinear interpolation



Feature  $f_{xy}$  for point (x, y)is a linear combination of features at its four neighboring grid cells:

(e.g. 3 x 640 x 480)

(e.g. 512 x 20 x 15)

He et al, "Mask R-CNN", ICCV 2017

$$f_{xy} = \sum_{i,j=1}^{2} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_j|)$$

Sample at regular points in each subregion using bilinear interpolation

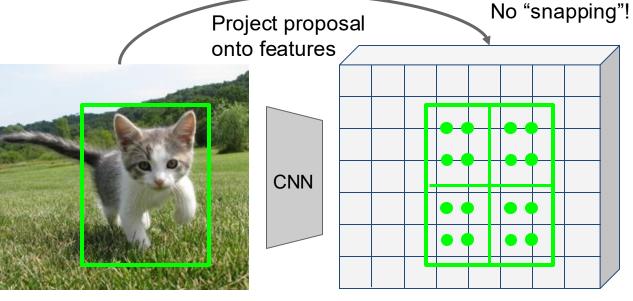
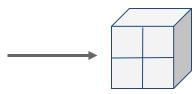


Image features: C x H x W (e.g. 512 x 20 x 15)

Max-pool within each subregion

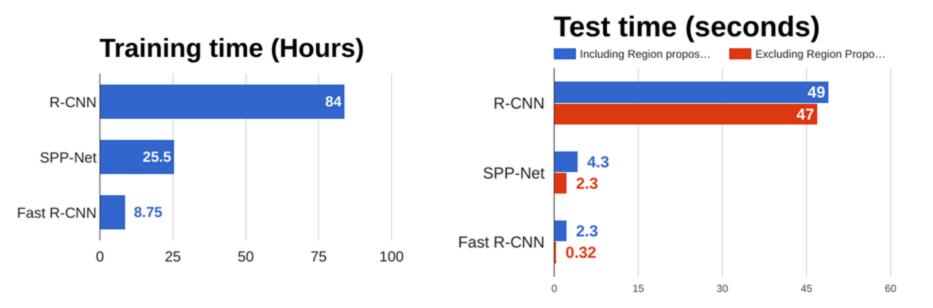


Region features (here 512 x 2 x 2; In practice e.g 512 x 7 x 7)

Input Image

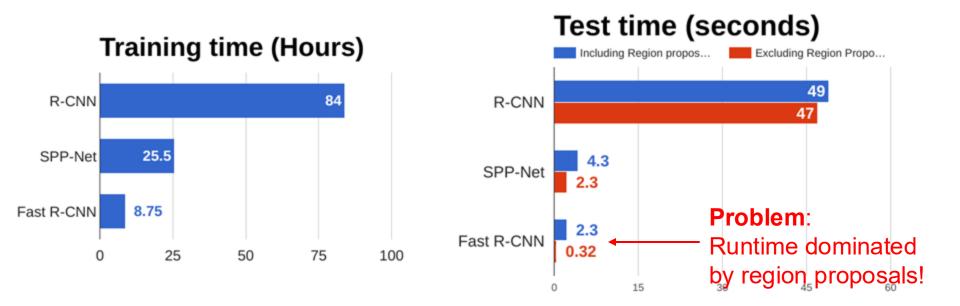
 $(e.g. 3 \times 640 \times 480)$ 

#### R-CNN vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014 Girshick, "Fast R-CNN", ICCV 2015

#### R-CNN vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014 Girshick, "Fast R-CNN", ICCV 2015

# Faster R-CNN: Make CNN do proposals!

Insert Region Proposal
Network (RPN) to predict
proposals from features

Otherwise same as Fast R-CNN: Crop features for each proposal, classify each one

Classification Bounding-box regression loss Classification **Bounding-box** Rol pooling loss regression loss proposals Region Proposal Network feature map CNN

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015 Figure copyright 2015, Ross Girshick; reproduced with permission

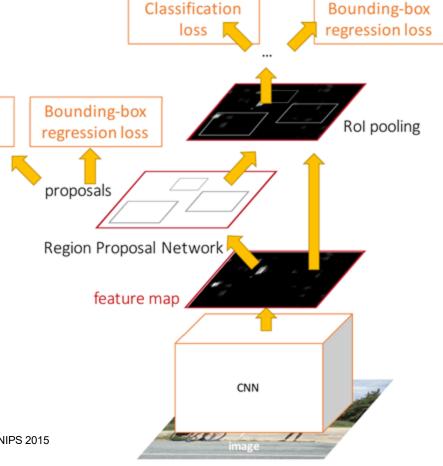
## Faster R-CNN: Make CNN do proposals!

Classification loss

Jointly train with 4 losses:

1. RPN classify object / not object

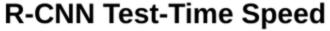
- 2. RPN regress box coordinate residuals
- 3. Final classification score (object classes)
- 4. Final box coordinates

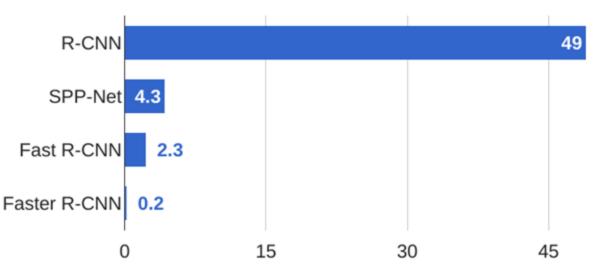


Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015 Figure copyright 2015, Ross Girshick; reproduced with permission

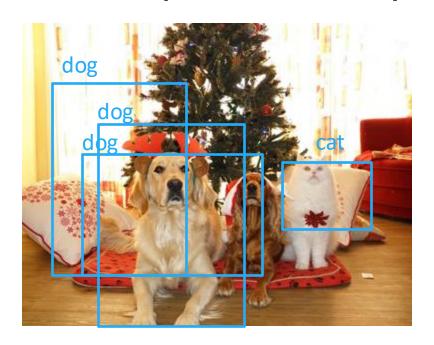
Faster R-CNN:

Make CNN do proposals!





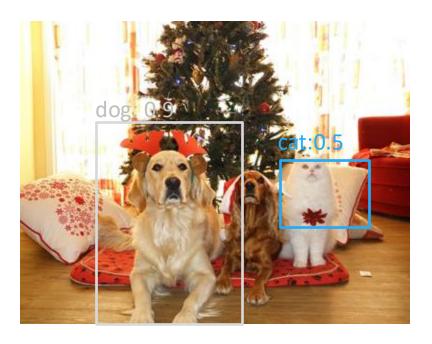
## Which prediction to pick?



**Problem**: Detectors almost always generate more box predictions than the number of objects in the image!
E.g., 300 is an upper bound how many objects we wish to detect.

We need to remove the **redundant predictions!** 

## Non-Max Suppression (NMS)



Intuitively: locally pick the box that has the highest "objectless" or class score and suppress other boxes that have significant overlap with the chosen box

Step 1: Pick highest-score prediction box Step 2: Remove bounding boxes with Intersection over Union (IoU) scores higher than certain threshold (e.g., 0.5)

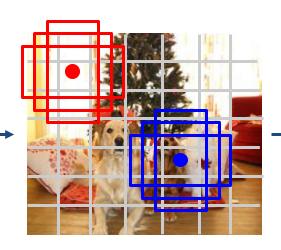
Go back to step 1

## Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Input image 3 x H x W

Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016 Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016 Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017



Divide image into grid 7 x 7

Image a set of **base boxes** centered at each grid cell Here B = 3

#### Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers:
   (dx, dy, dh, dw, confidence)
- Predict scores for each of C classes (including background as a class)
- Looks a lot like RPN, but category-specific!

Output: 7 x 7 x (5 \* B + C)

### Object Detection: Lots of variables ...

Backbone Network

VGG16

ResNet-101

Inception V2

Inception V3

Inception

ResNet

MobileNet

"Meta-Architecture"

Two-stage: Faster R-CNN Single-stage: YOLO / SSD

Hybrid: R-FCN

Image Size # Region Proposals

- - -

**Takeaways** 

Faster R-CNN is slower but more accurate

SSD is much faster but not as accurate

Bigger / Deeper backbones work better

Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017 Zou et al, "Object Detection in 20 Years: A Survey", arXiv 2019

R-FCN: Dai et al, "R-FCN: Object Detection via Region-based Fully Convolutional Networks", NIPS 2016
Inception-V2: Ioffe and Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shiff", ICML 2015
Inception V3: Szegedy et al, "Rethinking the Inception Architecture for Computer Vision", arXiv 2016
Inception ResNet: Szegedy et al, "Inception-V4, Inception-ResNet and the Impact of Residual Connections on Learning", arXiv 2016
MobileNet: Howard et al, "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

## Instance Segmentation

Classification



CAT

No spatial extent

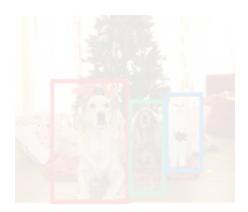
Semantic Segmentation



TREE, SKY

No objects, just pixels

**Object Detection** 



DOG, DOG, CAT

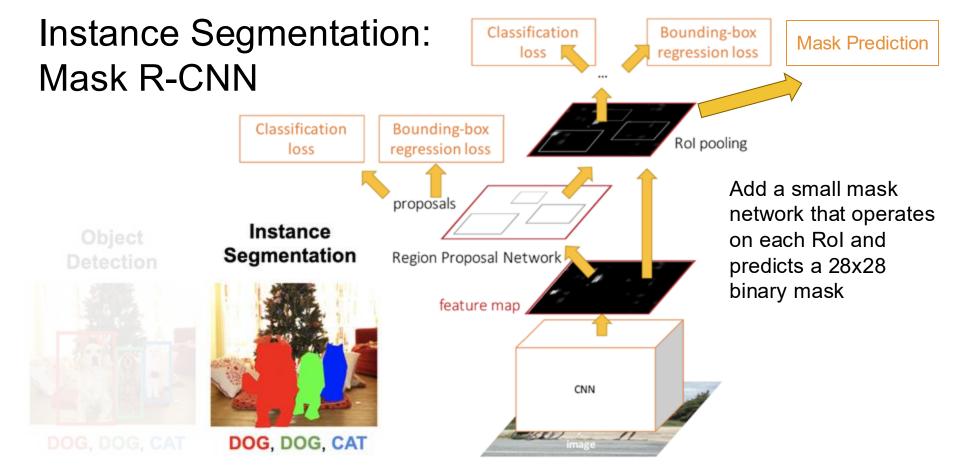
Instance Segmentation



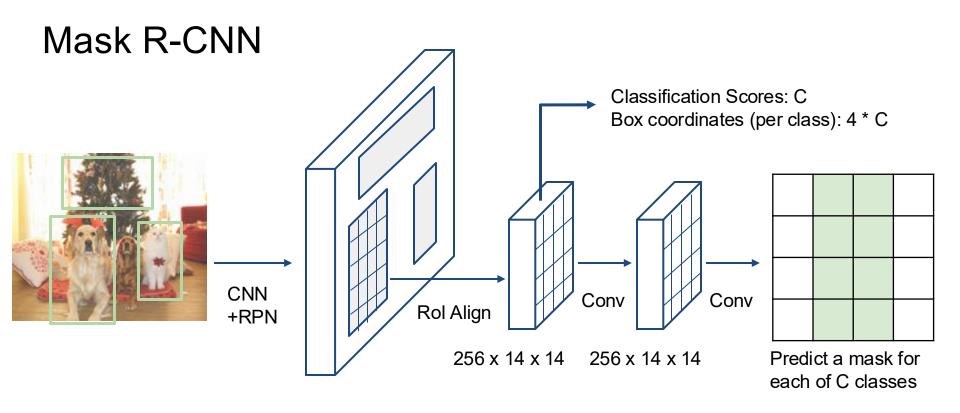
DOG, DOG, CAT

Multiple Object

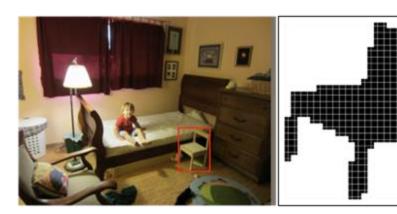
Object Detection: Classification Bounding-box regression loss loss **Faster R-CNN** Classification Bounding-box Rol pooling regression loss loss proposals Object Region Proposal Network Detection feature map CNN DOG, DOG, CAT



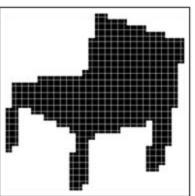
He et al, "Mask R-CNN", ICCV 2017



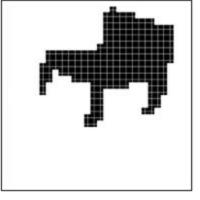
C x 28 x 28



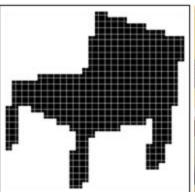




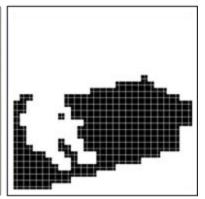








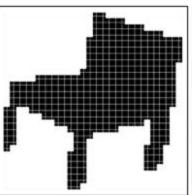




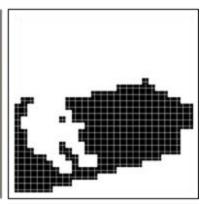








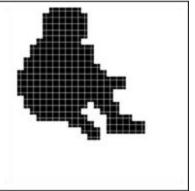






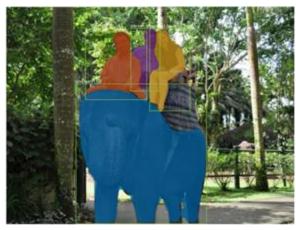






## Mask R-CNN: Very Good Results!



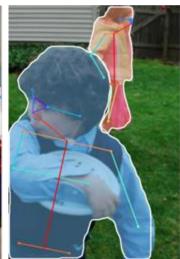




## Mask R-CNN Also does pose





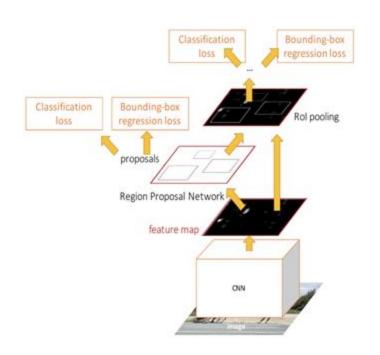


#### **RCNN Series**

- R-CNN: Per-region detection, handcrafted region proposal
- Fast R-CNN: Shared feature extraction, Rol Pooling, Anchors
- Faster R-CNN: Region Proposal Networks, Rol Align
- Mask R-CNN: Instance Segmentation

Detectors are becoming more complex! Many hyperparameters to tune for each components ...

Can we simplify it?



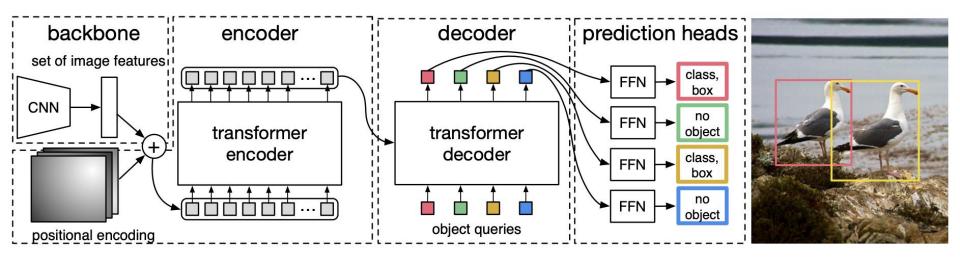
#### End-to-End Object Detection with Transformers

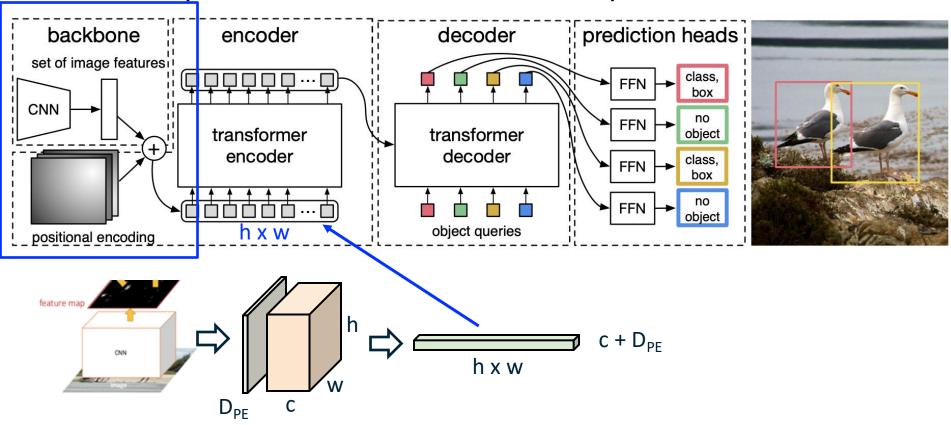
Nicolas Carion\*, Francisco Massa\*, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko

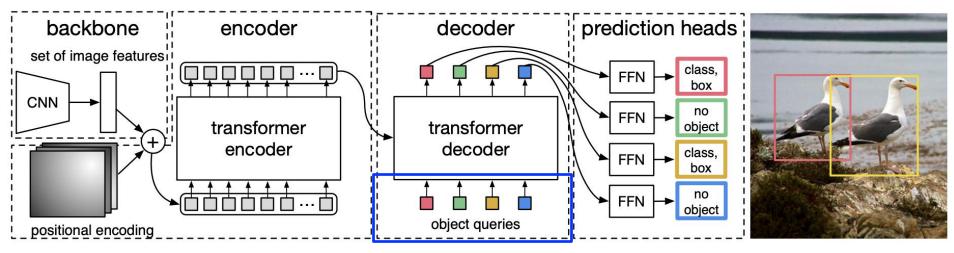
Facebook AI

#### **Key ideas:**

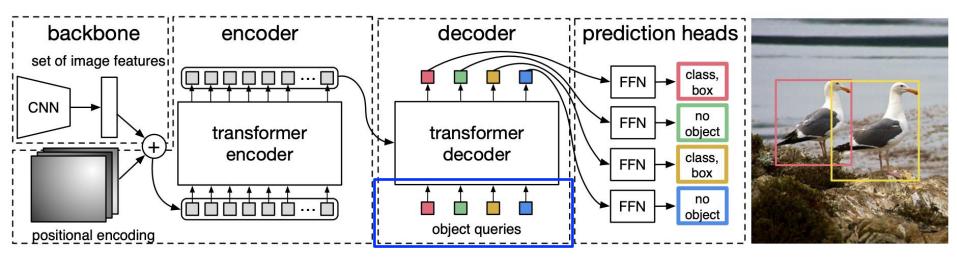
- Detection as a set-to-set prediction problem
- Use Transformer to model the detection problem

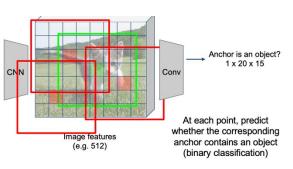






A fixed set of *learnable embeddings*, e.g., 300 size-N vectors Q: Why?



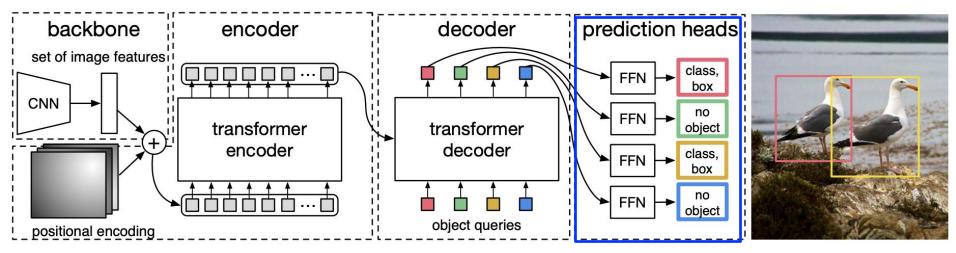


A fixed set of learnable embeddings, e.g., 300 size-N vectors

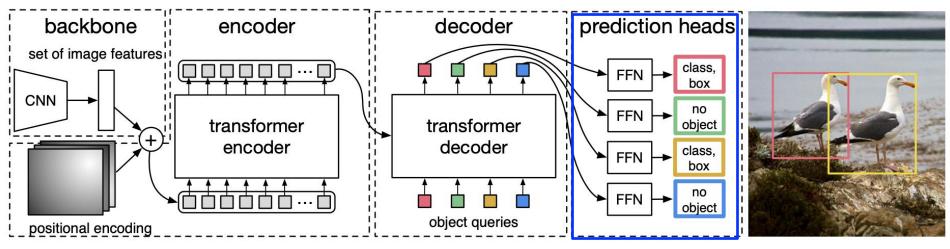
Q: Why?

A: Break the symmetry of predictions, so that each prediction is different.

Analogous to anchors in \*R-CNN, but no spatial location



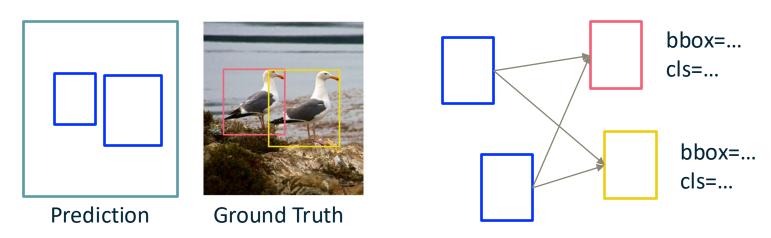
**Problem**: We don't know which query corresponds to which ground truth during training! We can't predetermine a fixed order like in sequence decoding.



**Problem**: We don't know which query corresponds to which ground truth during training! We can't predetermine a fixed order like in sequence decoding.

**Solution:** Set matching loss --- train your model to generate a set of predictions that matches ground truth regardless of its order.

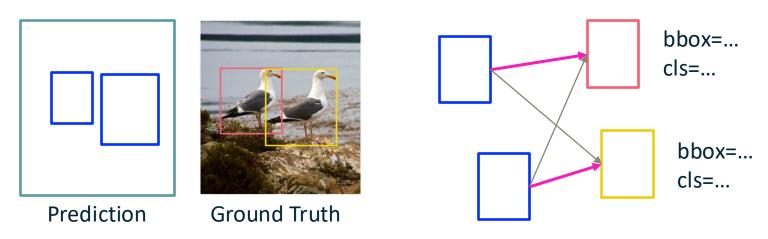
## Hangarian Loss (Set Matching Loss)



Goal: minimize bipartite distance

**Problem**: each query should be trained to match one ground truth. We don't know the matching!

## Hangarian Loss (Set Matching Loss)



Goal: minimize bipartite distance

- 1. **Hungarian matching:** find the minimum-loss bipartite matching between prediction and ground truth given the current prediction.
- 2. **Minimize matched loss:** Given the matched prediction and ground truth, minimize the detection loss (bounding box distance and classification CE loss)

#### DETR vs. FasterRCNN

Model	GFLOPS/FPS	#params	AP	$\mathrm{AP}_{50}$	AP <sub>75</sub>	$\mathrm{AP_S}$	$AP_{M}$	$\overline{\mathrm{AP_L}}$
Faster RCNN-DC5	320/16	166M	39.0	60.5	42.3	21.4	43.5	52.5
Faster RCNN-FPN	180/26	42M	40.2	61.0	43.8	24.2	43.5	52.0
Faster RCNN-R101-FPN	246/20	60M	42.0	62.5	45.9	25.2	45.6	54.6
Faster RCNN-DC5+	320/16	166M	41.1	61.4	44.3	22.9	45.9	55.0
Faster RCNN-FPN+	180/26	42M	42.0	62.1	45.5	26.6	45.4	53.4
Faster RCNN-R101-FPN+	246/20	60M	44.0	63.9	<b>47.8</b>	27.2	48.1	56.0
DETR	86/28	41M	42.0	62.4	44.2	20.5	45.8	61.1
DETR-DC5	187/12	41M	43.3	63.1	45.9	22.5	47.3	61.1
DETR-R101	152/20	60M	43.5	63.8	46.4	21.9	48.0	61.8
DETR-DC5-R101	253/10	60M	44.9	64.7	47.7	23.7	49.5	<b>62.3</b>

Similar size, simpler, and (mostly) better!

We are still detecting a fixed number of object with finite vocabulary ...

#### **Segment Anything**

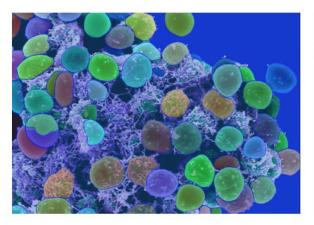
Alexander Kirillov<sup>1,2,4</sup> Eric Mintun<sup>2</sup> Nikhila Ravi<sup>1,2</sup> Hanzi Mao<sup>2</sup> Chloe Rolland<sup>3</sup> Laura Gustafson<sup>3</sup> Tete Xiao<sup>3</sup> Spencer Whitehead Alexander C. Berg Wan-Yen Lo Piotr Dollár<sup>4</sup> Ross Girshick<sup>4</sup>

<sup>1</sup>project lead <sup>2</sup>joint first author <sup>3</sup>equal contribution <sup>4</sup>directional lead

Meta AI Research, FAIR

#### **Key ideas:**

- Query-based prediction instead of fixed set-to-set prediction
- Large-scale training data with auto-labeling







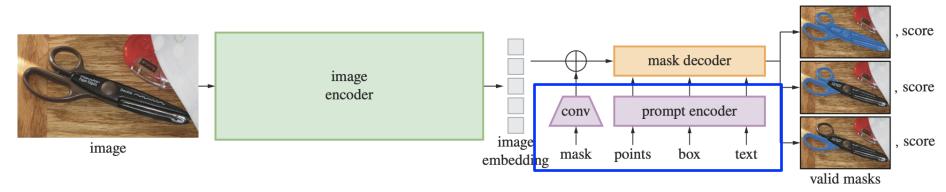
SegmentAnything (Meta AI, 2023)

Try it yourself! https://segment-anything.com/demo#



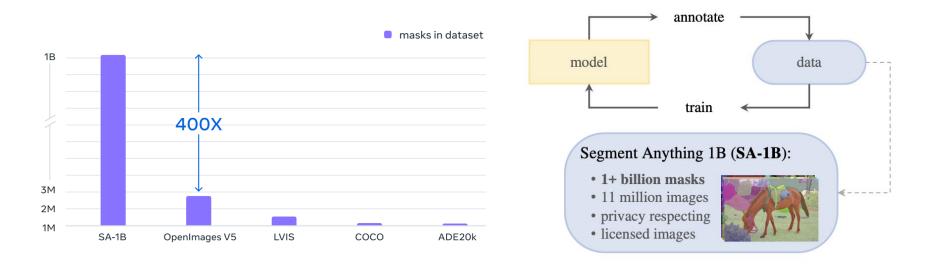
SegmentAnything (Meta AI, 2023)

Try it yourself! https://segment-anything.com/demo#



No more learned embeddings. Query anything you want!

SegmentAnything (Meta AI, 2023)



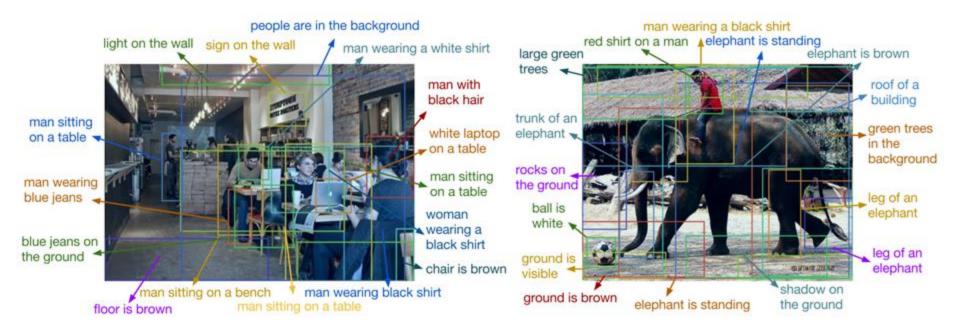
SegmentAnything (Meta AI, 2023)

#### Summary: Segmentation and Detection

- Object segmentation and detection are most common application of computer vision research.
- They have driven decades of advancement in Autonomous Vehicles, Robotics, traffic analytics, and basically any devices that have camera and adequate computing power (e.g., Smart Phone)
- Segmentation and Detection with DNNs evolved through similar paths (sliding window, feature sharing, input-output, alignment etc.)
- We have a new wave of foundation detection and segmentation models driven by Transformer + ConvNet + large dataset

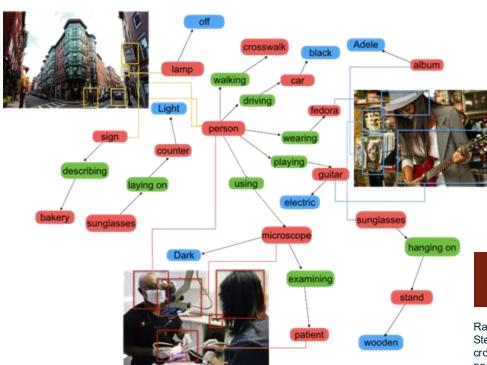
# Beyond 2D Object Detection...

# Object Detection + Captioning = Dense Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016 Figure copyright IEEE, 2016. Reproduced for educational purposes.

## Objects + Relationships = 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

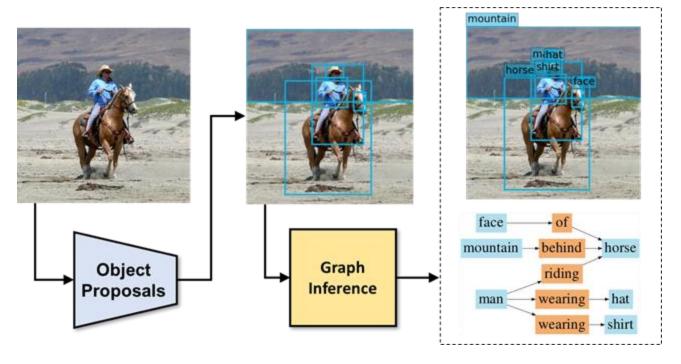
2.3 Million Relationships

Everything Mapped to Wordnet Synsets

# **VISUAL**GENOME

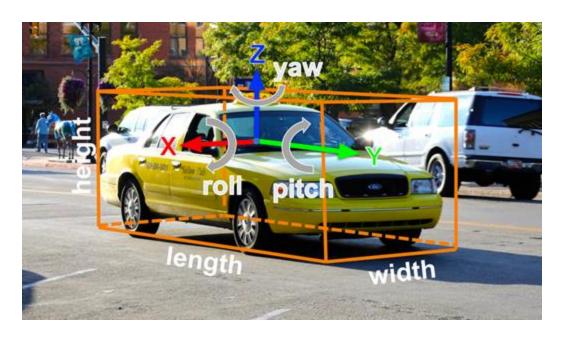
Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations." International Journal of Computer Vision 123, no. 1 (2017): 32-73.

#### Scene Graph Prediction



Xu, Zhu, Choy, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", CVPR 2017 Figure copyright IEEE, 2018. Reproduced for educational purposes.

#### 3D Object Detection



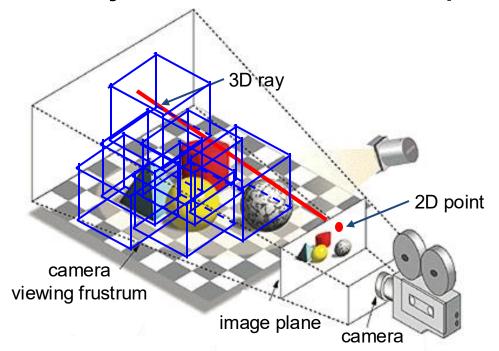
2D Object Detection: 2D bounding box (x, y, w, h)

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: Simple Camera Model

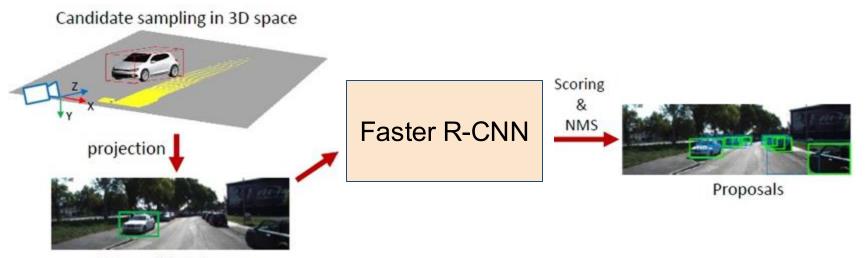


A point on the image plane corresponds to a ray in the 3D space

A 2D bounding box on an image is a **frustrum** in the 3D space

Localize an object in 3D: The object can be anywhere in the **camera viewing frustrum!** 

#### 3D Object Detection: Monocular Camera



#### 2D candidate boxes

- 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.

#### How to Represent 3D Data?

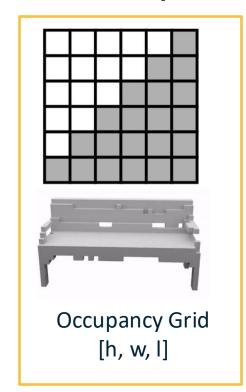


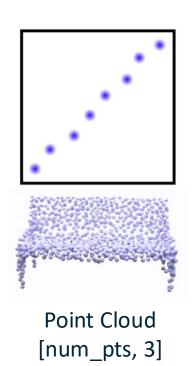


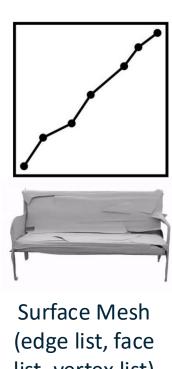


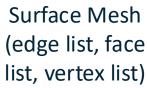
7

#### 3D Representations









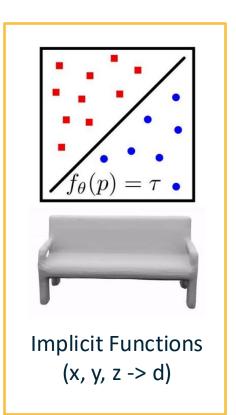
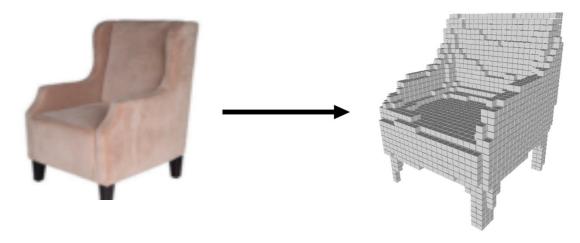


Figure credit: Autonomous Vision Group

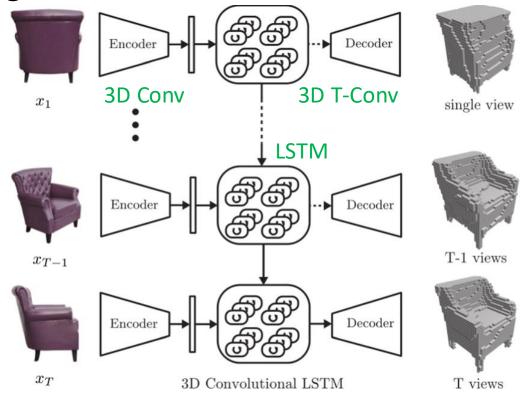
#### 3D Occupancy Grid



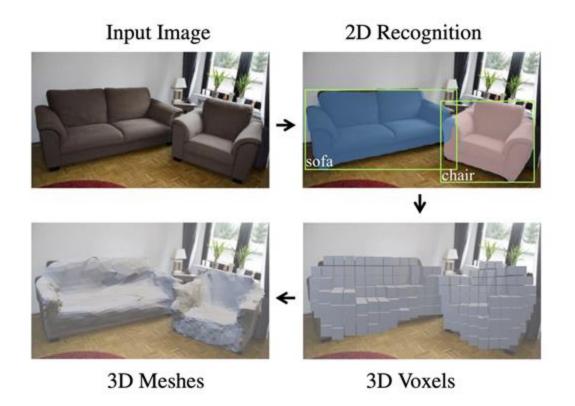
Represent the "occupancy" of objects in 3D space with a 3D voxel grid

- $V \in \{0, 1\}^{[H, W, L]}$
- Just like segmentation in Masked-RCNN, but in 3D!
- Conceptually simple
- Not trivial to scale to high-resolution shapes

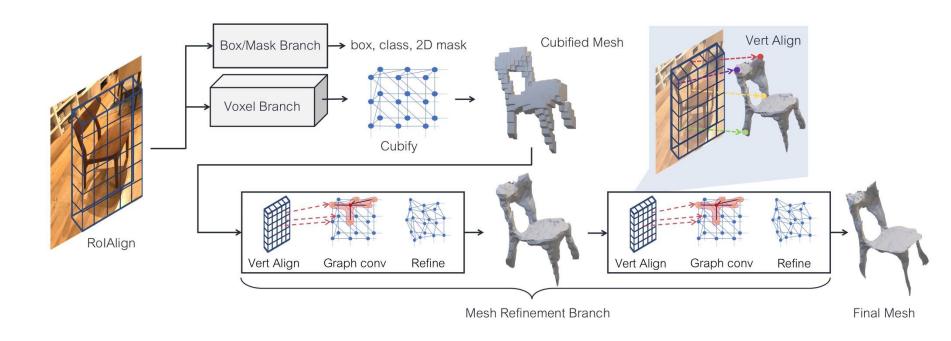
#### Predicting 3D Voxel Grid with 3D ConvNet



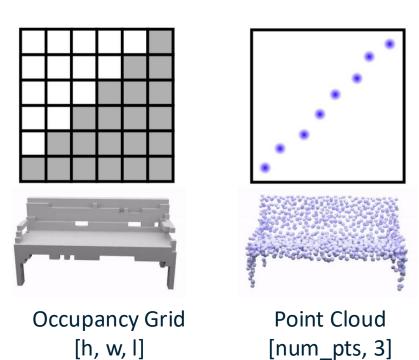
#### Detection + Reconstruction: Mesh R-CNN



#### Detection + Reconstruction: Mesh R-CNN



#### 3D Representations



Surface Mesh (edge list, face list, vertex list)

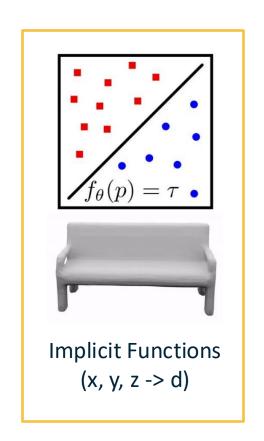
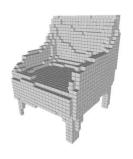


Figure credit: Justin Johnson

Example: representing a 3D occupancy grid

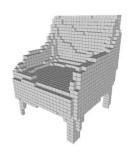


**Explicit**: A tensor of **3D voxel grid**  $V \in \{0, 1\}^{[H,W,L]}$ 

**Implicit:** A **function** that maps locations to occupancies

 $F_{\theta}$ :  $x, y, z \rightarrow \{0, 1\}$ 

Example: representing a 3D occupancy grid



**Explicit**: A tensor of **3D voxel grid**  $V \in \{0, 1\}^{[H,W,L]}$ 

**Implicit:** A **function** that maps locations to occupancies

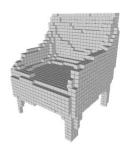
 $F_{\theta}$ :  $x, y, z \rightarrow \{0, 1\}$ 

Implicit representation describes 3D shapes using **mathematical functions** rather than explicit voxels, points, or mesh.

Example: Signed Distance Function

 $F_{\theta} \colon \mathbb{R}^3 \to \mathbb{R}$ 

Example: representing a 3D occupancy grid



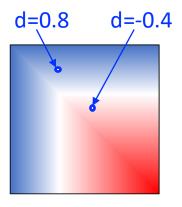
**Explicit**: A tensor of **3D voxel grid**  $V \in \{0, 1\}^{[H,W,L]}$ 

**Implicit:** A **function** that maps locations to occupancies

$$F_{\theta}$$
:  $x, y, z \rightarrow \{0, 1\}$ 

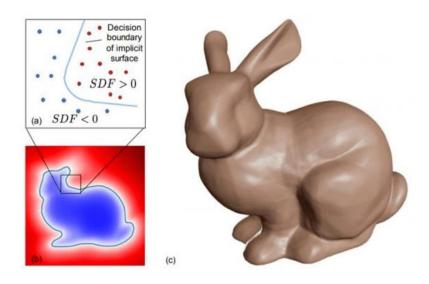
Implicit representation describes 3D shapes using **mathematical functions** rather than explicit voxels, points, or mesh.

Example: Signed Distance Function  $F_{\theta} \colon \mathbb{R}^{N} \to \mathbb{R}$ 



How far is a point from the nearest surface, and is the point *inside or outside* of the shape?

SDF distance map



DeepSDF representation applied to the Stanford Bunny: (a) depiction of the underlying implicit surface SDF = 0 trained on sampled points inside SDF < 0 and outside SDF > 0 the surface, (b) 2D cross-section of the signed distance field, (c) rendered 3D surface recovered from SDF = 0.

Can we train NNs to represent more than just geometry of a 3D shape?

Park et al., DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation

#### Implicit 3D Representation: Beyond Geometry

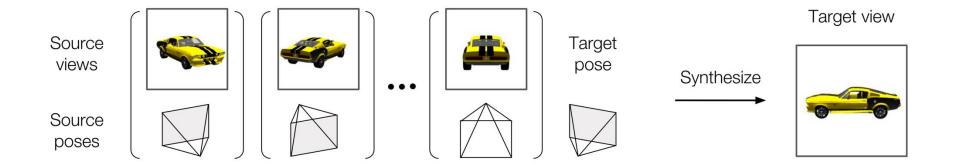


 $f_{\theta}(viewpoint) = Image$ 

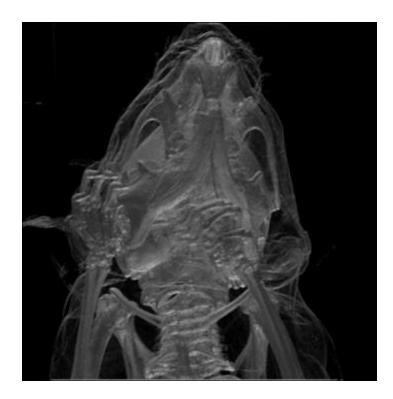
**Goal**: Learn an implicit 3D representation function that maps any camera viewpoint to full RGB images

Can we implicitly represent a full 3D scene, including its fine-grained geometry (e.g., surface occupancy) and appearance?

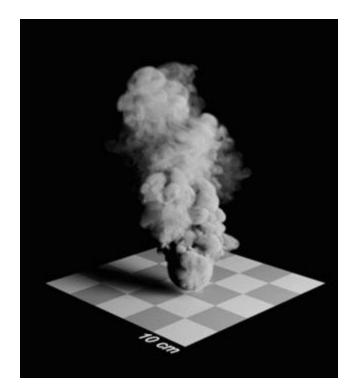
## Problem: Novel View Synthesis



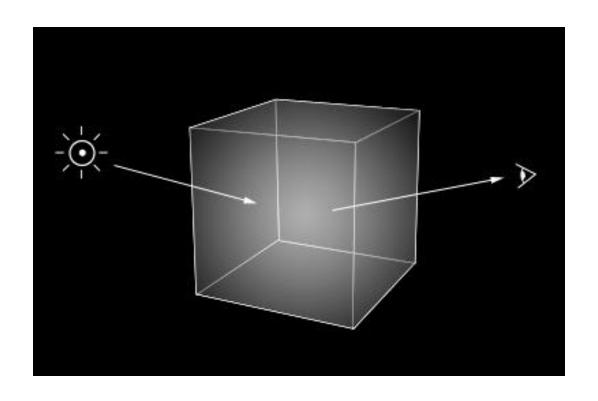
### **Basics: Volume Rendering**

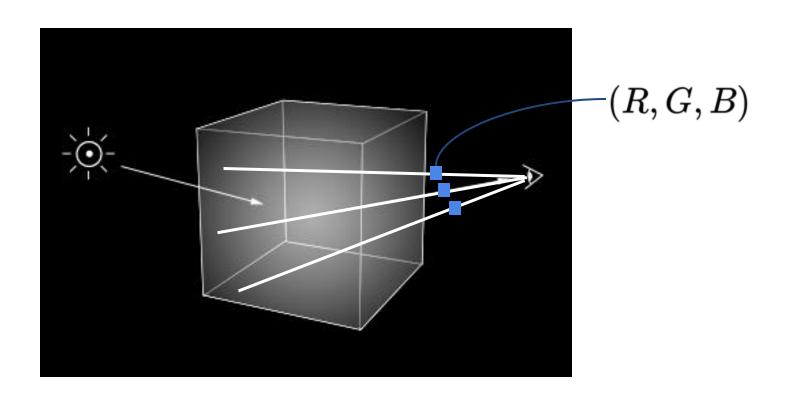


https://en.wikipedia.org/wiki/Volume\_rendering

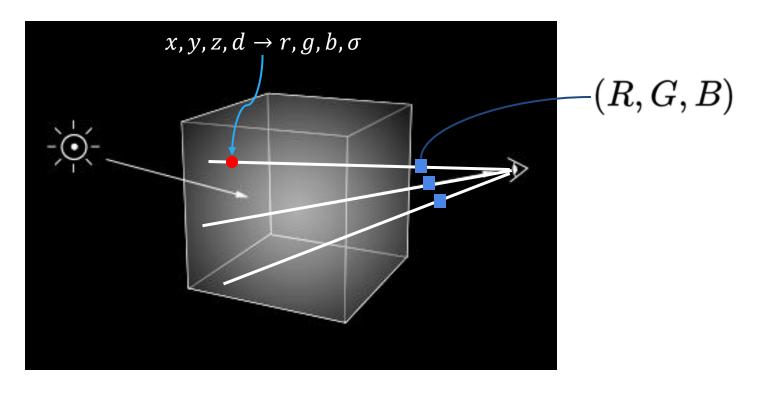


https://coronarenderer.freshdesk.com/support/solutions/articles/12000045276-how-to-use-the-corona-volume-grid-

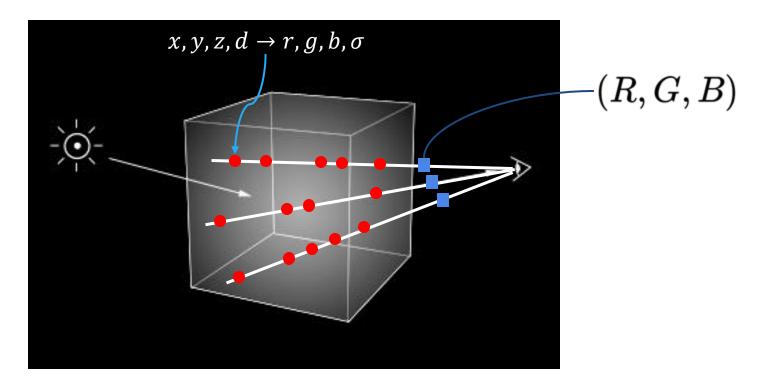




Each location (x, y, z) emits certain color r, g, b when viewed with direction d. We represent point occupancy continuously as density d.

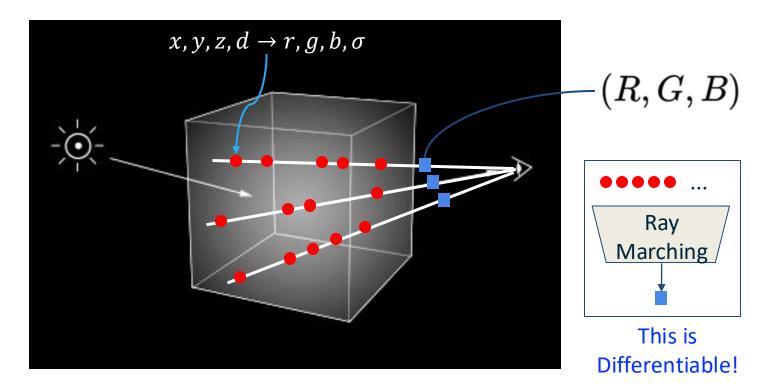


Each location (x, y, z) emits certain color r, g, b when viewed with direction d. We represent point occupancy continuously as density d.



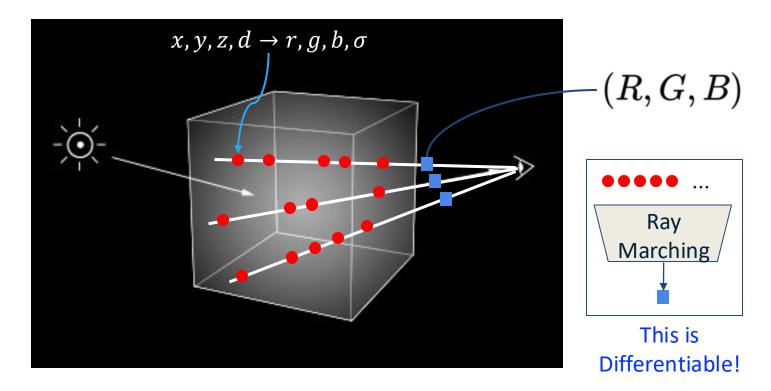
#### Volume Rendering: Ray Marching

**Ray Marching**: Integrate color and density of points along a ray (via discretization) to render an RGB value. Render many points -> An image!



#### Volume Rendering: Ray Marching

**Idea:** Train a neural network to represent the ray marching volume rendering function:  $F_{\theta}(x, y, z, d) \rightarrow (r, g, b, \sigma)$ . **Each NN encodes a 3D scene**.

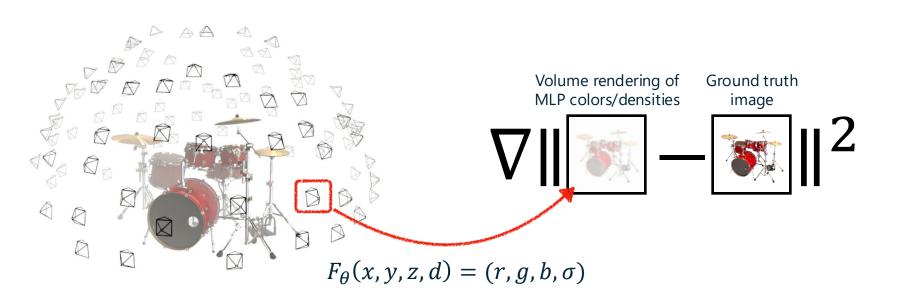


#### NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis

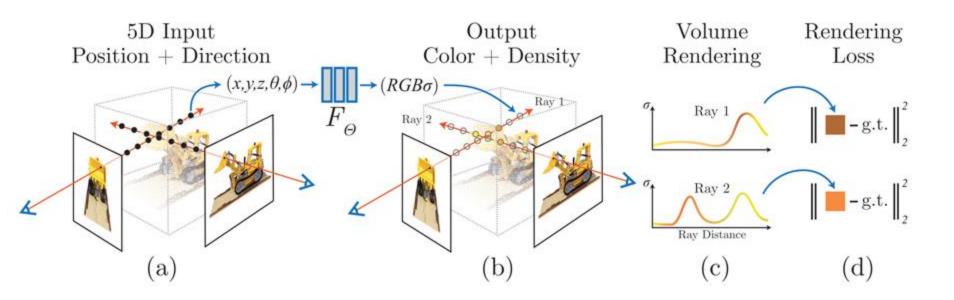
Ben Mildenhall<sup>1\*</sup> Pratul P. Srinivasan<sup>1\*</sup> Matthew Tancik<sup>1\*</sup> Jonathan T. Barron<sup>2</sup> Ravi Ramamoorthi<sup>3</sup> Ren Ng<sup>1</sup>

<sup>1</sup>UC Berkeley <sup>2</sup>Google Research <sup>3</sup>UC San Diego

# Train a Single Neural Network to Reproduce the Ground Truth Images of a Scene



#### **NeRF** Overview



#### **NeRF: Optimization**

The volume density  $\sigma(\mathbf{x})$  can be interpreted as the differential probability of a ray terminating at an infinitesimal particle at location  $\mathbf{x}$ . The expected color  $C(\mathbf{r})$  of camera ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  with near and far bounds  $t_n$  and  $t_f$  is:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t),\mathbf{d})dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s))ds\right).$$
(1)

Solution: Numerically estimate the integral (quadrature).

- 1. Discretize the ray into bins.
- 2. Sample point in each bin.
- 3. Compute numerical integration.

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^{N} T_i (1 - \exp(-\sigma_i \delta_i)) c_i, \text{ where } T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right)$$

### Key Insight 1: Positional Encoding

Challenge: Having  $F_{\theta}$  operate directly on (x, y, z, d) performs poorly.

Solution: Positional encoding

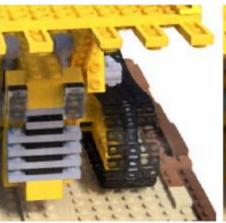
$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p))$$



Ground Truth



Complete Model





No View Dependence No Positional Encoding

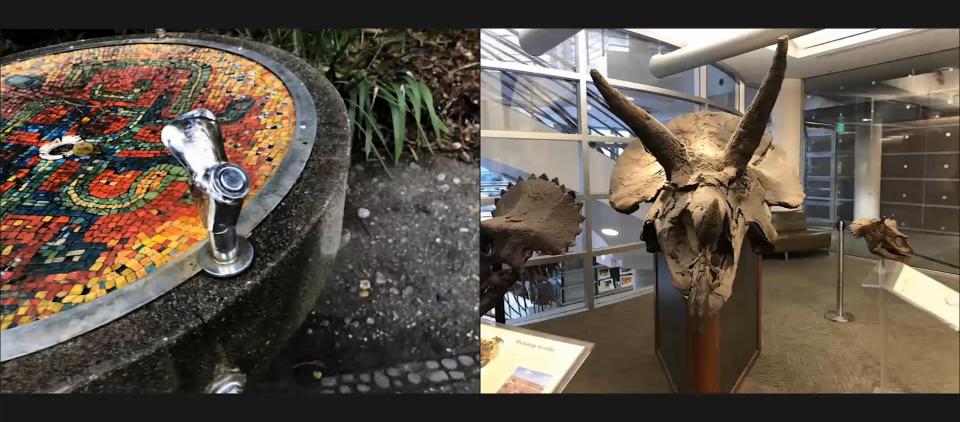
### Key Insight 2: Hierarchical Volume Rendering

Challenge: Waste of compute on empty space.

Solution: coarse-to-fine prediction.

$$\hat{C}_c(\mathbf{r}) = \sum_{i=1}^{N_c} w_i c_i, \quad w_i = T_i (1 - \exp(-\sigma_i \delta_i)).$$
 (5)

Normalizing these weights as  $\hat{w}_i = w_i/\sum_{j=1}^{N_c} w_j$  produces a piecewise-constant PDF along the ray. We sample a second set of  $N_f$  locations from this distribution using inverse transform sampling, evaluate our "fine" network at the union of the first and second set of samples, and compute the final rendered color of the ray  $\hat{C}_f(\mathbf{r})$  using Eqn. 3 but using all  $N_c + N_f$  samples. This procedure allocates more



# NeRF encodes convincing view-dependent effects using directional dependence



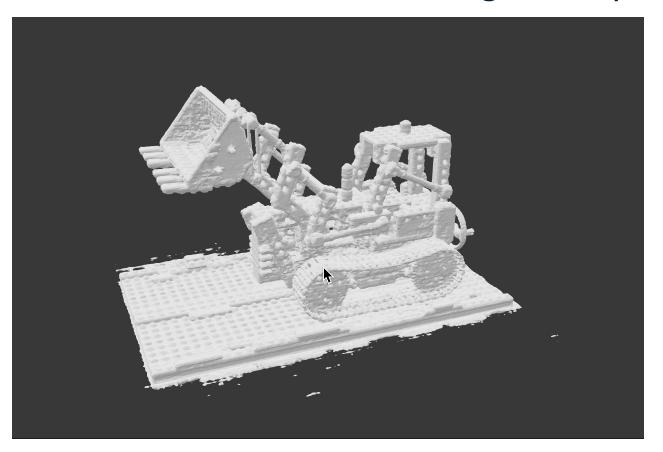
# NeRF encodes convincing view-dependent effects using directional dependence



#### NeRF encodes detailed scene geometry with occlusion effects



#### NeRF encodes detailed scene geometry



#### 3D Gaussian Splatting (Kerbl and Kopanas et al., 2023)

Key idea: 3D Gaussians as an **explicit representation** of a scene

- Train Gaussian blobs via inverse rendering (similar to NeRF)
- Store scene as Gaussian blobs instead of neural network weights (NeRF)
- Much faster during inference, but takes a lot of space to store

#### NeRF Gaussian Splatting







#### Summary: 3D Representation and Neural Rendering

- Representation matters a lot for 3D computer vision tasks (detection, reconstruction, etc.)
- 3D Voxels are intuitive representation of space but struggles with highresolution shape and large scenes
- Implicit function emerge as a new paradigm in representing scenes with Neural Networks
- Neural volume rendering: represent scenes implicit as point-direction to color-density neural networks. Photorealistic rendering, slow to train and evaluate
- More recent works on trading off space and time

## Next Time: Vision Language Models