

facebook  
Artificial Intelligence

# Computer Vision @ Scale

Manohar Paluri  
Director, Facebook AI



A globe of the Earth is centered in the image, set against a dark, starry space background. The globe is overlaid with a complex network of thin, light-colored lines and small, multi-colored square nodes (in shades of gold, green, and red). The network is dense and interconnected, covering the entire surface of the globe and extending slightly beyond its edges, symbolizing global connectivity and digital networks.

Connecting People to bring them closer together







And helping machines understand the visual world is an important component!—

Before | **Text Input**



Now | **Visual Input**





What will you learn today?

How do you design a image and video recognition system for billion scale?

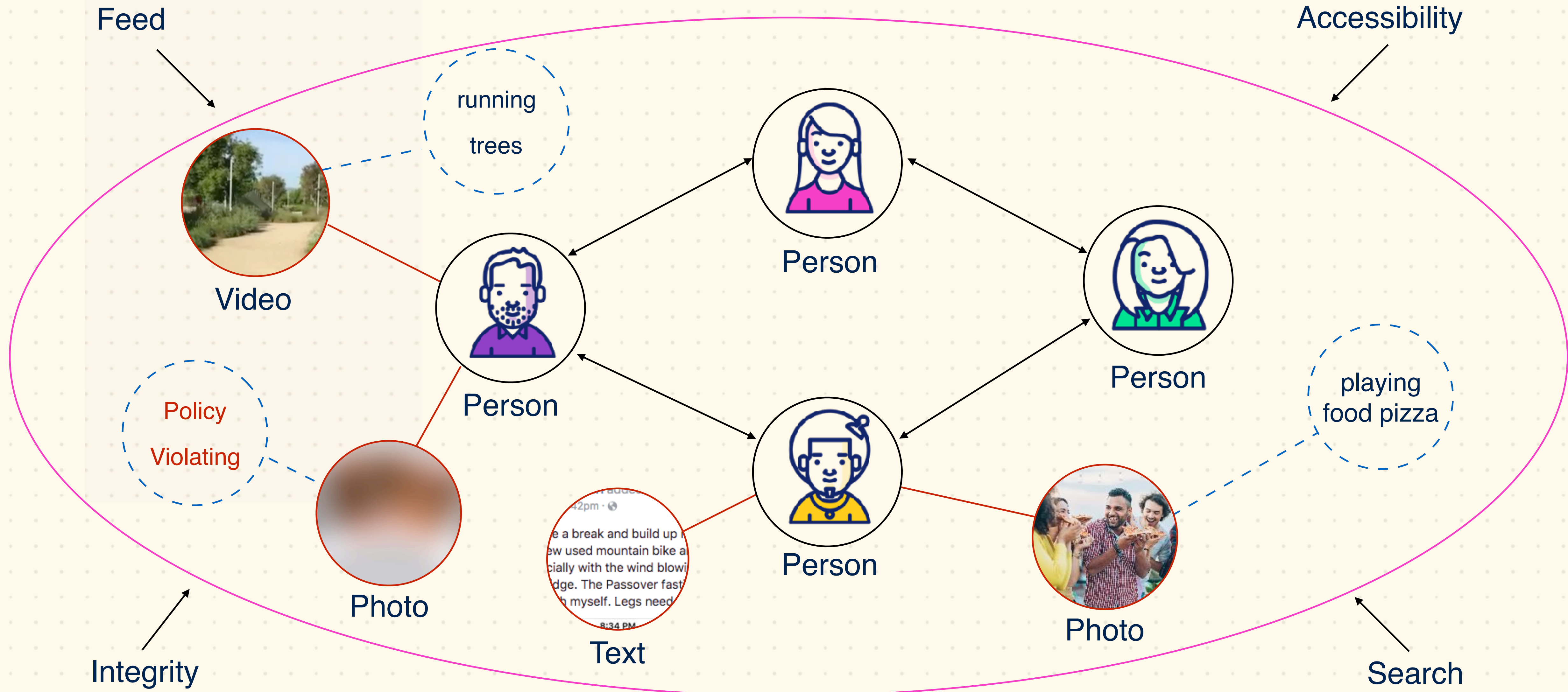
Can you remove the requirement of annotation to learn best representations?

Can we understand video faster than understanding individual frames?

How does pushing state of the art in CV make a meaningful difference to everyone in the world?



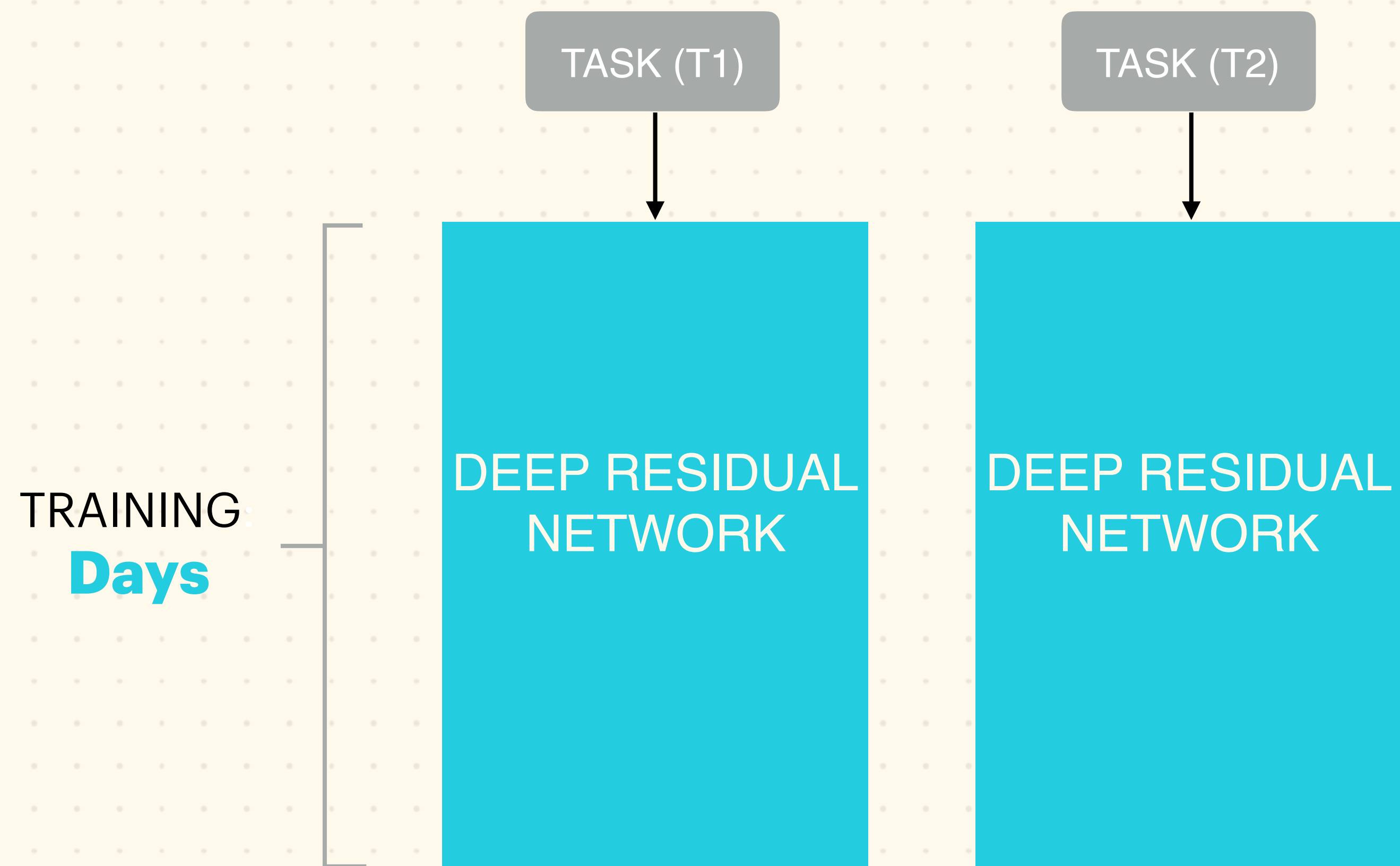
# Social Graph -> Semantic Graph





# Vision Models In Production

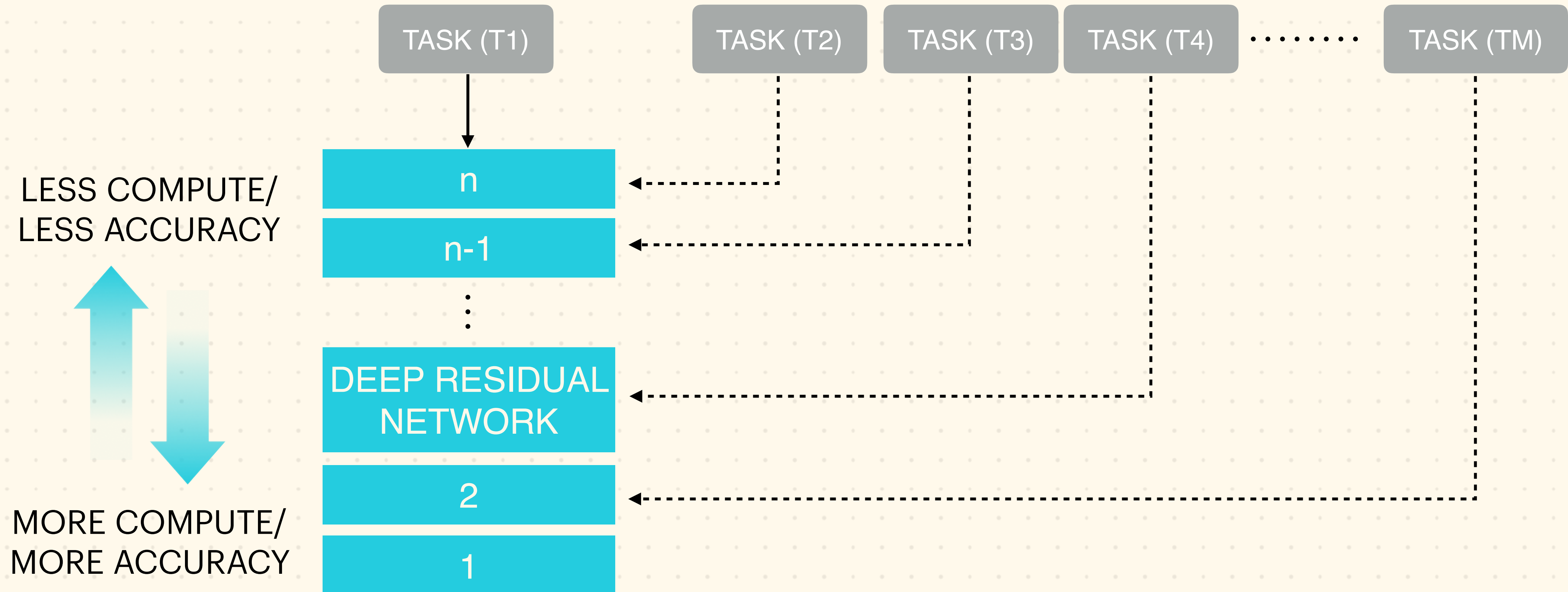
- Multiple vision tasks need to be done
- Cannot afford one separate model for each task
- Explosion in computations cost and resources





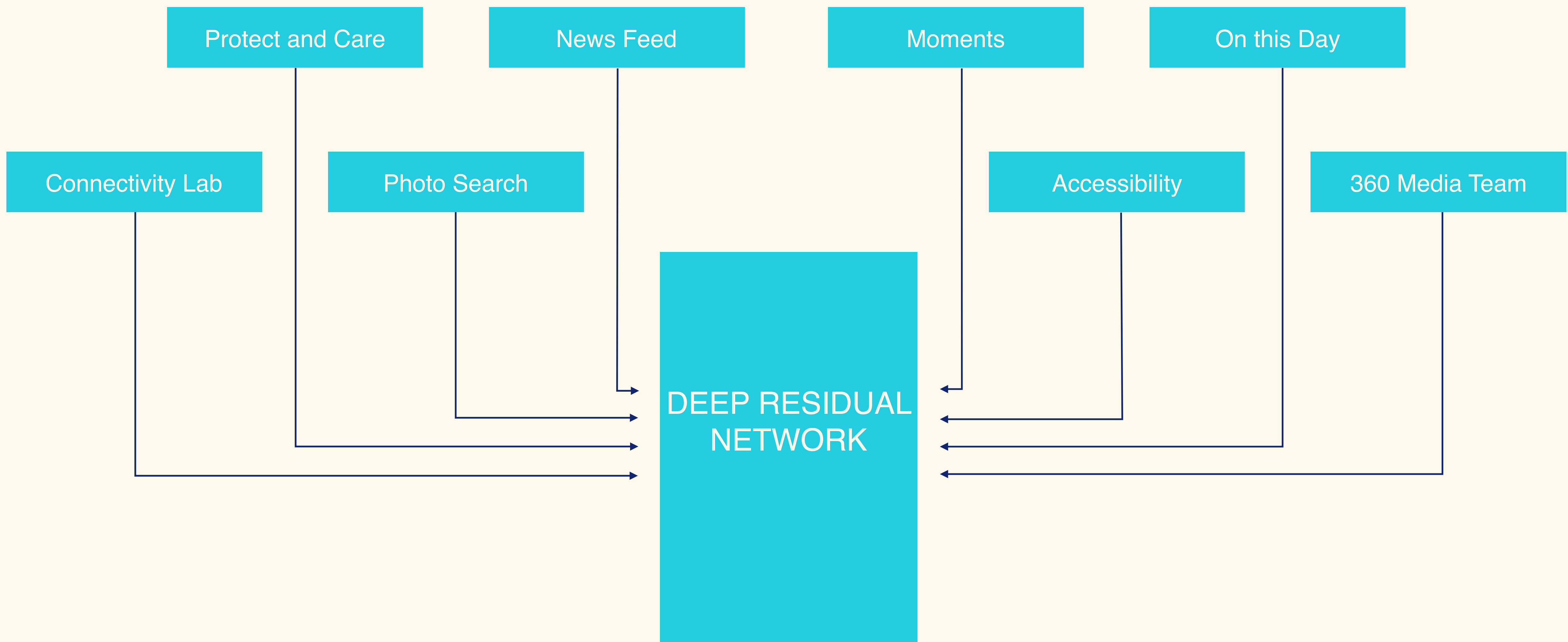
# Our Vision - Towards Universal Vision Model

SUBTITLE



Our work allows to move the tasks towards upper layers





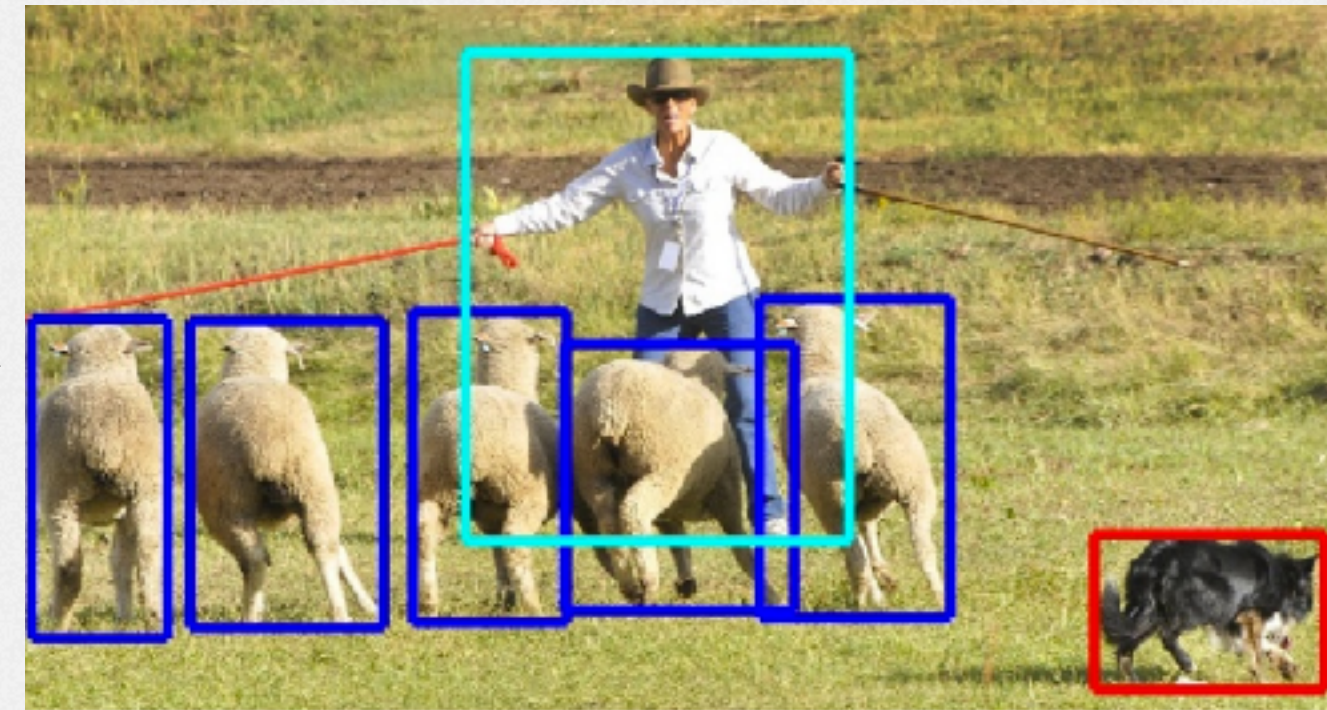
# Branches



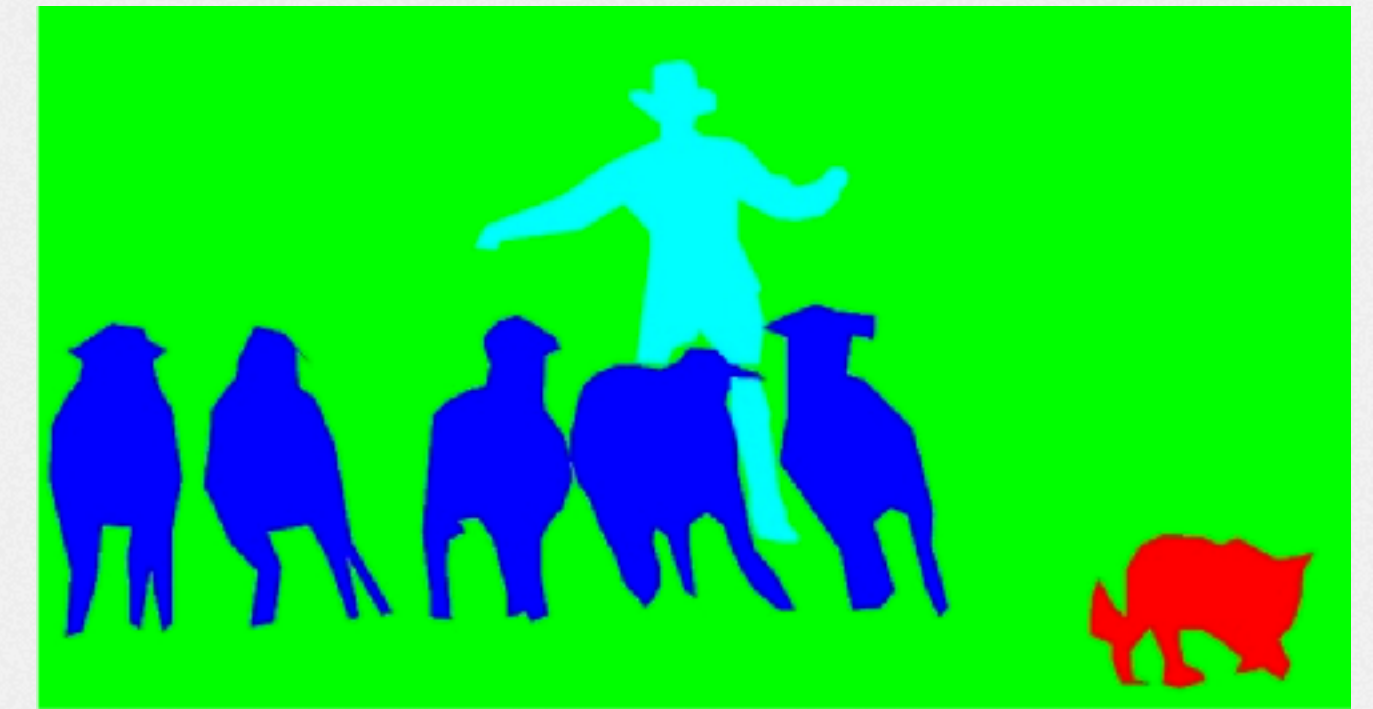
# Progress in Image Understanding



Classification



Detection



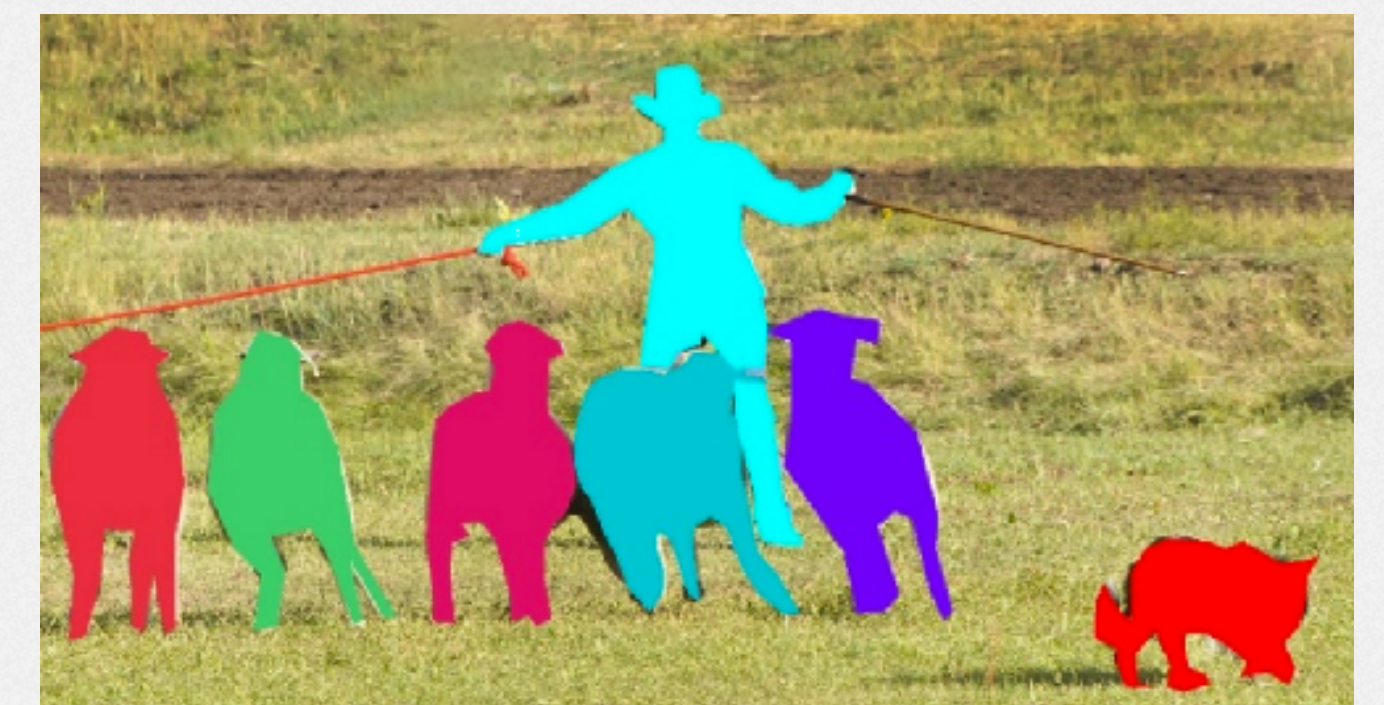
Semantic Segmentation



Relationships  
Fine grained  
Aesthetics



Attributes



Instance Segmentation

·  
·





ALEXNET

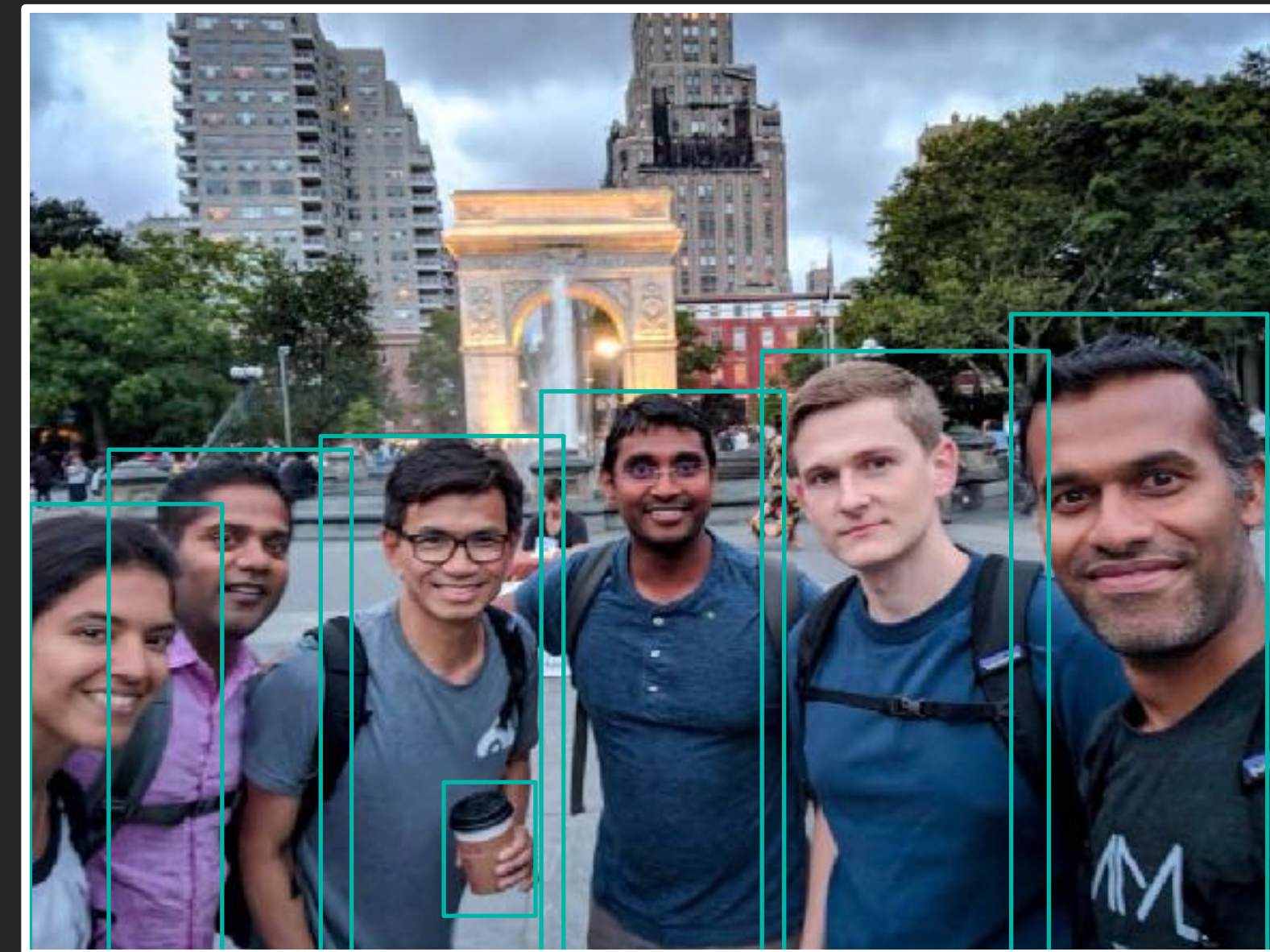
2012





ALEXNET

2012



FASTER R-CNN

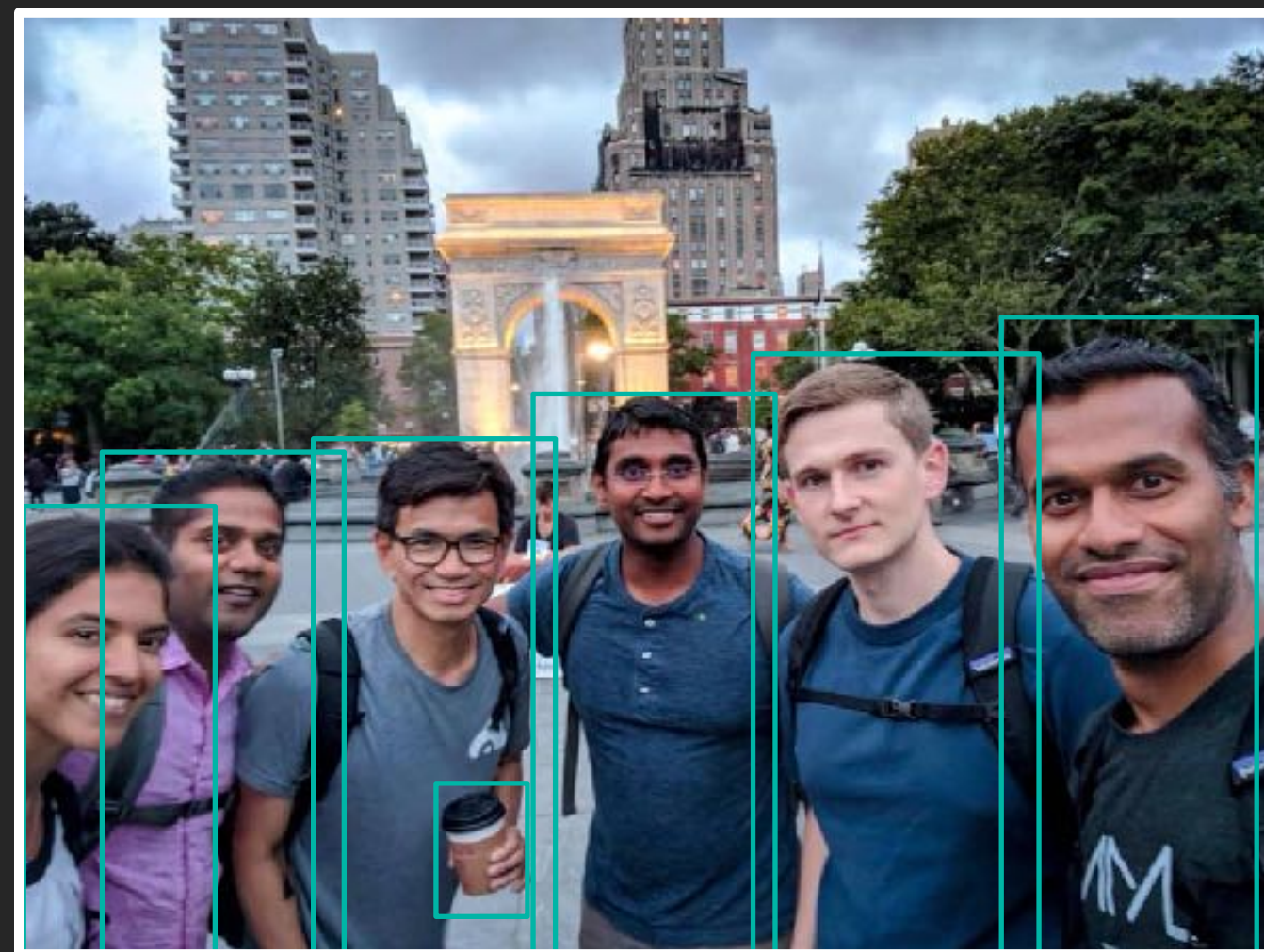
2015





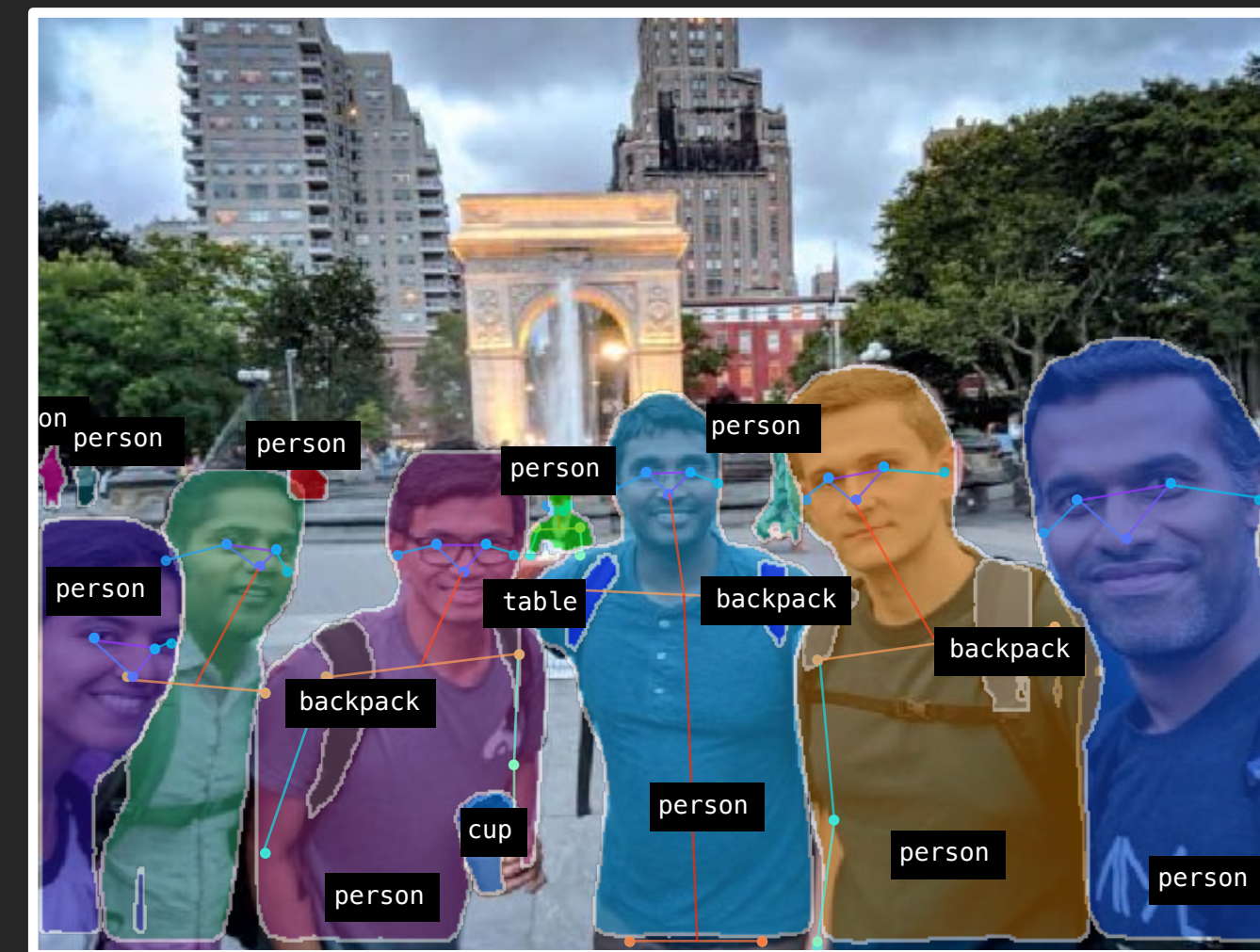
ALEXNET

2012



FASTER R-CNN

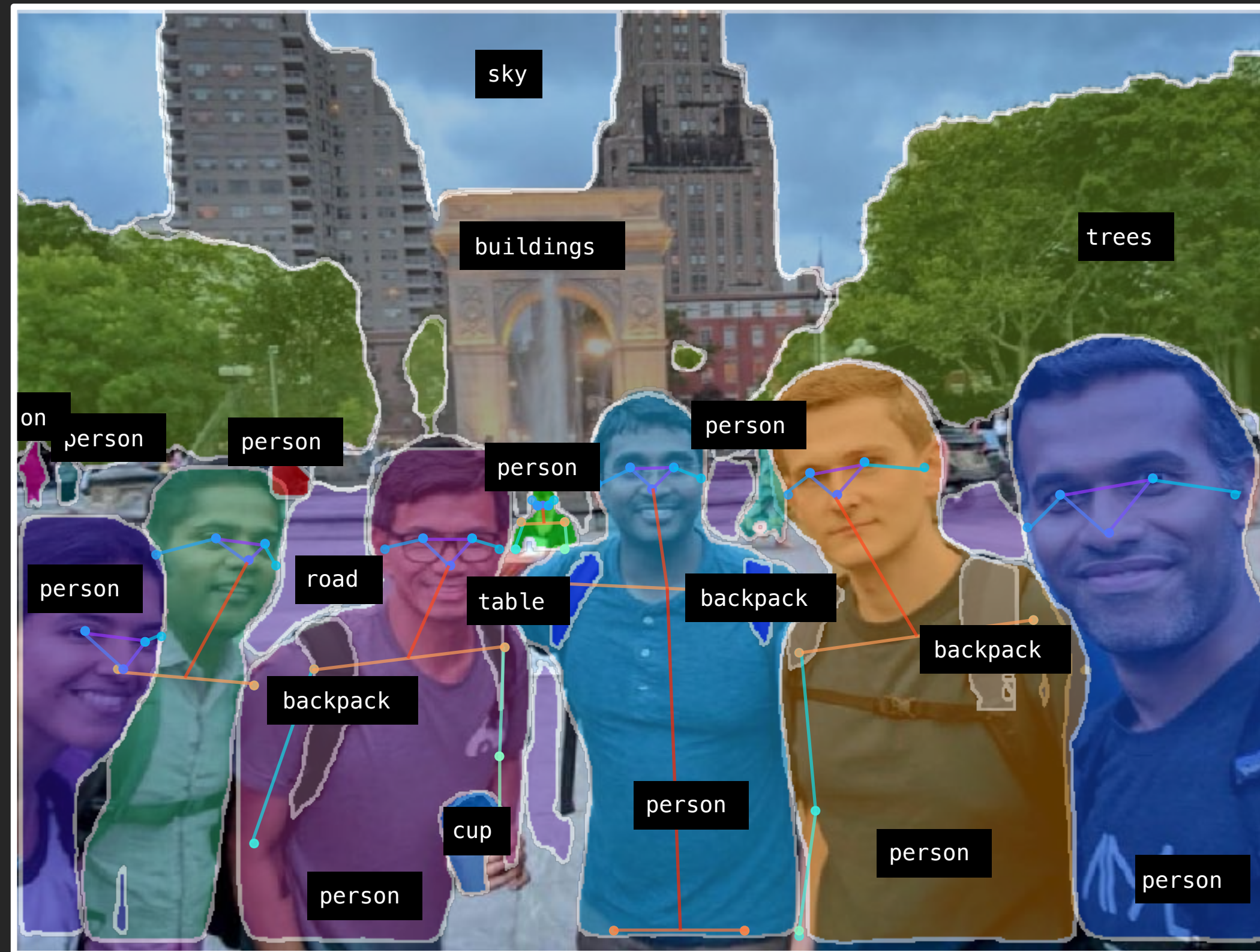
2015



MASK R-CNN

2017



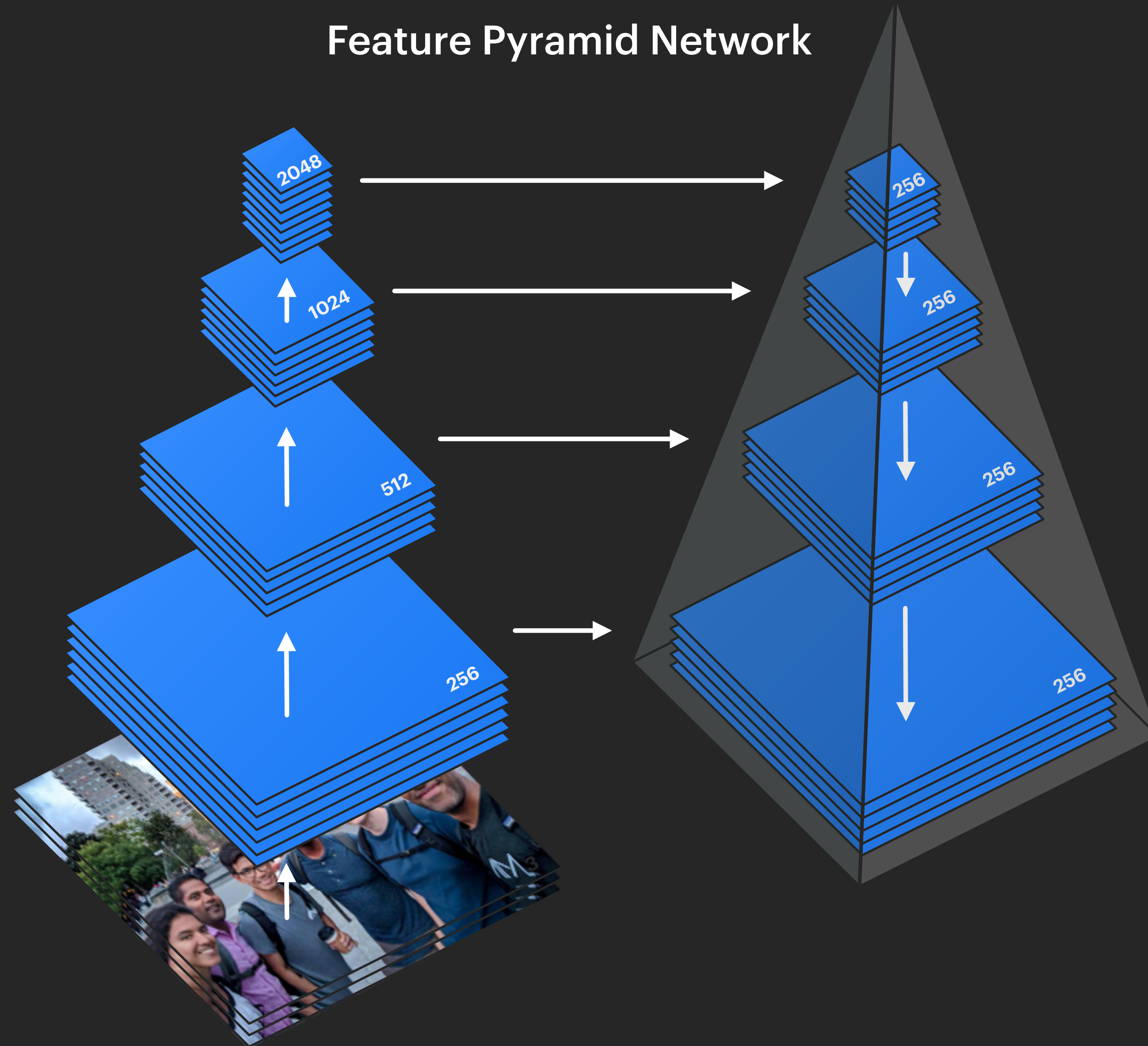


## PANOPTIC FPN

2019

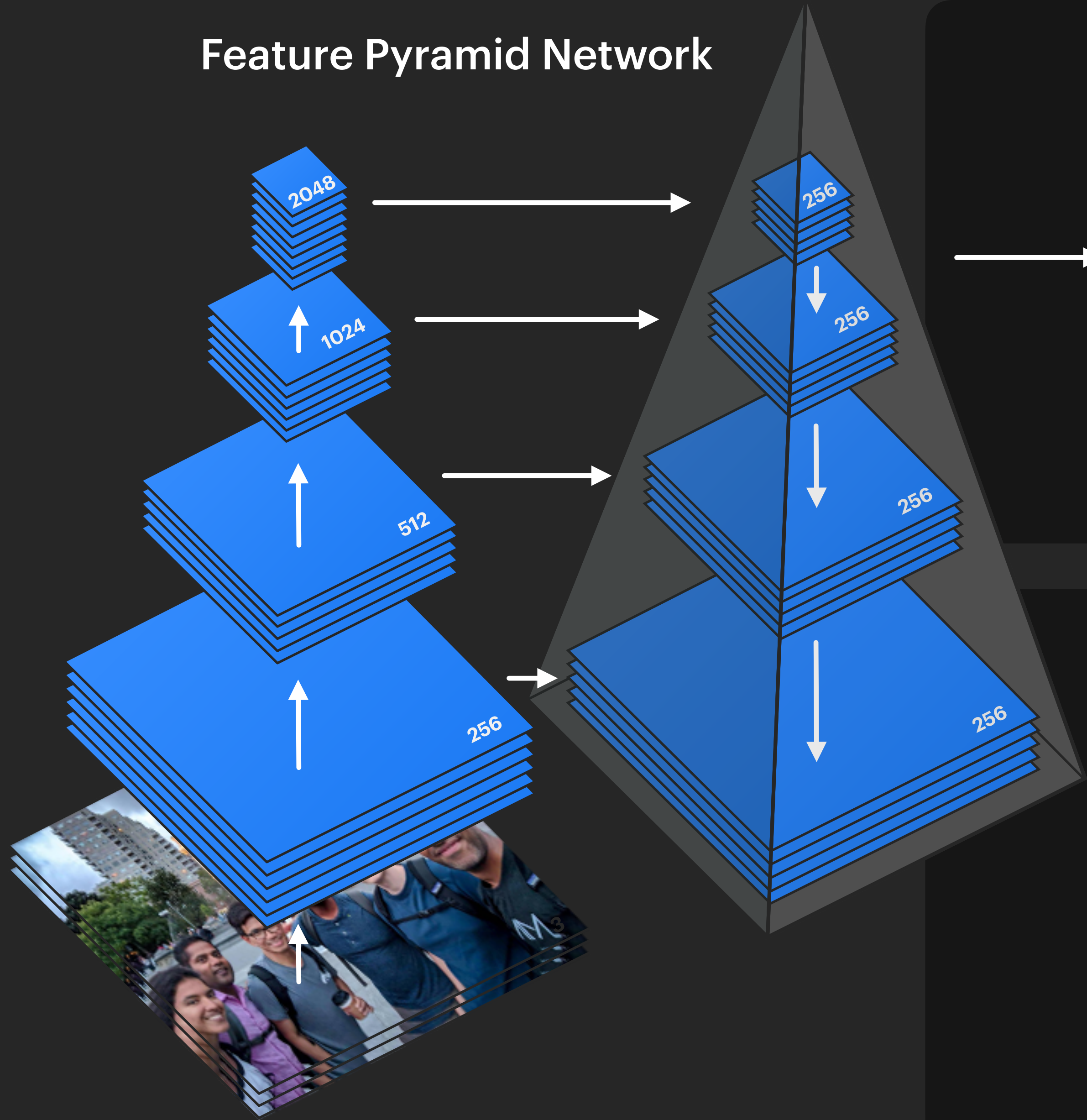


# Feature Pyramid Network

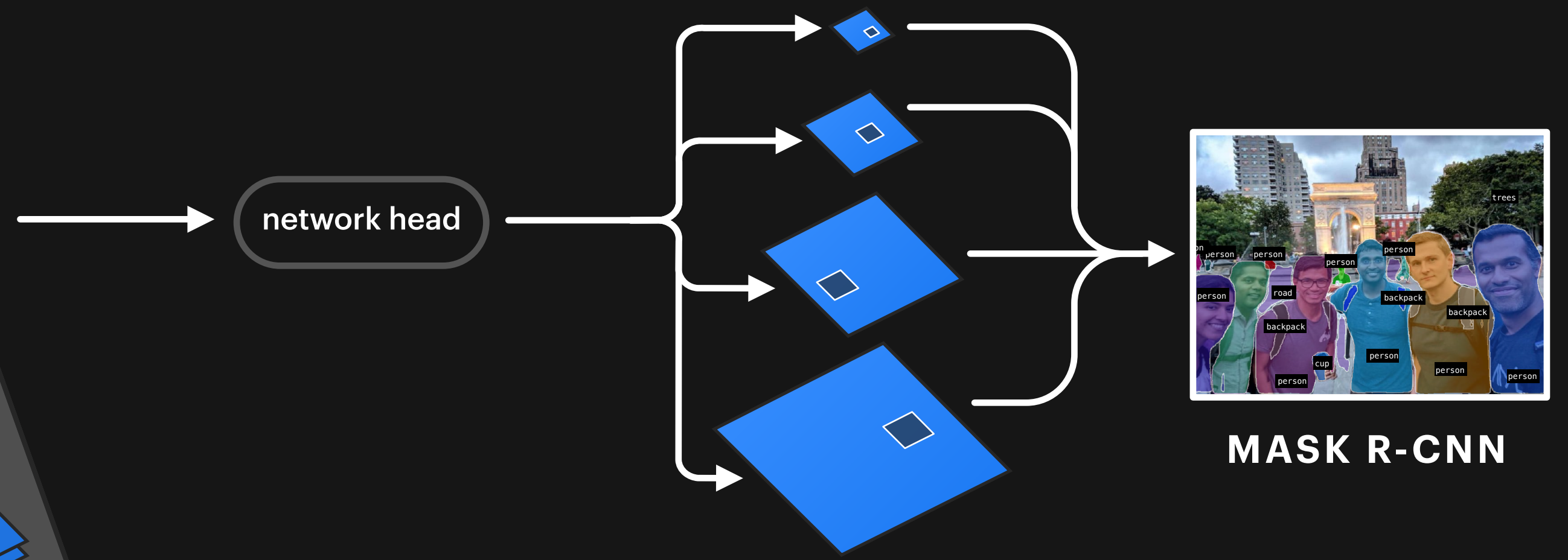




# Feature Pyramid Network

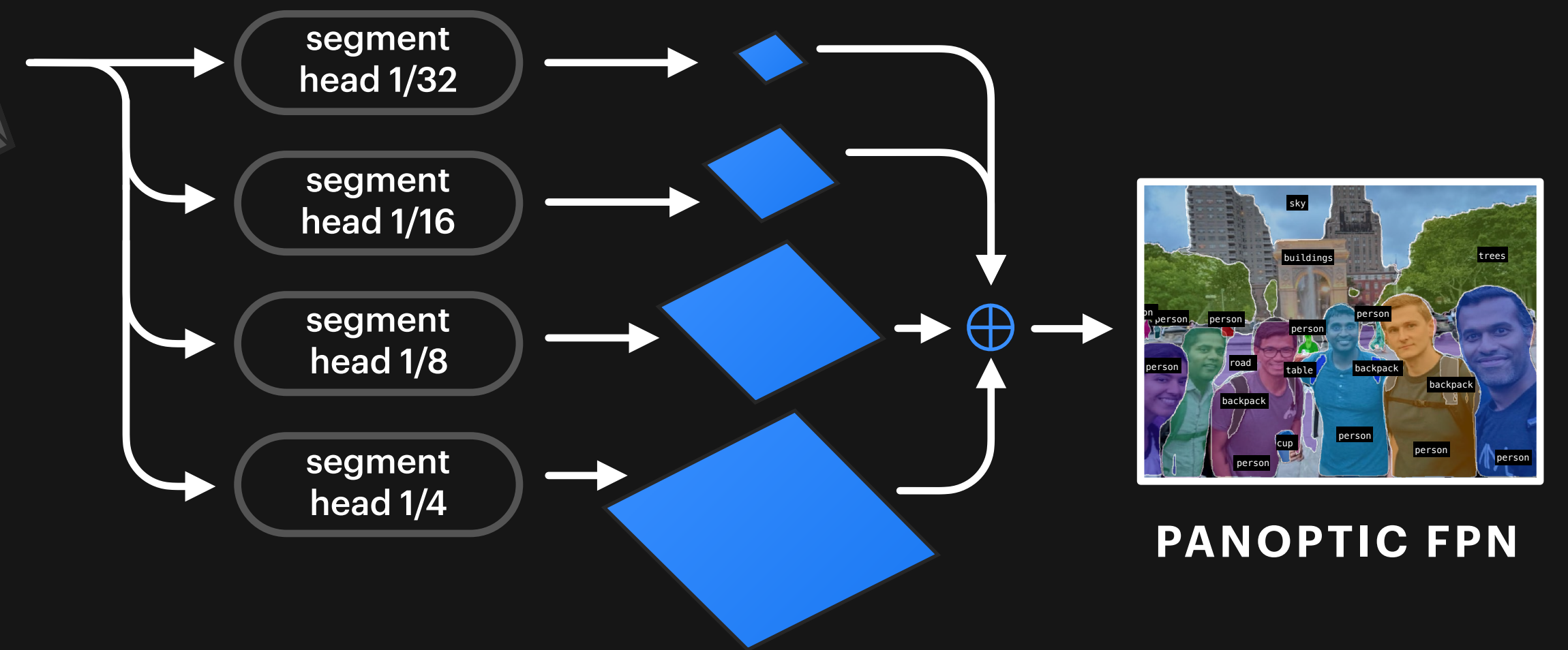


# Instance Segmentation Branch



**MASK R-CNN**

# Semantic Segmentation Branch



**PANOPTIC FPN**

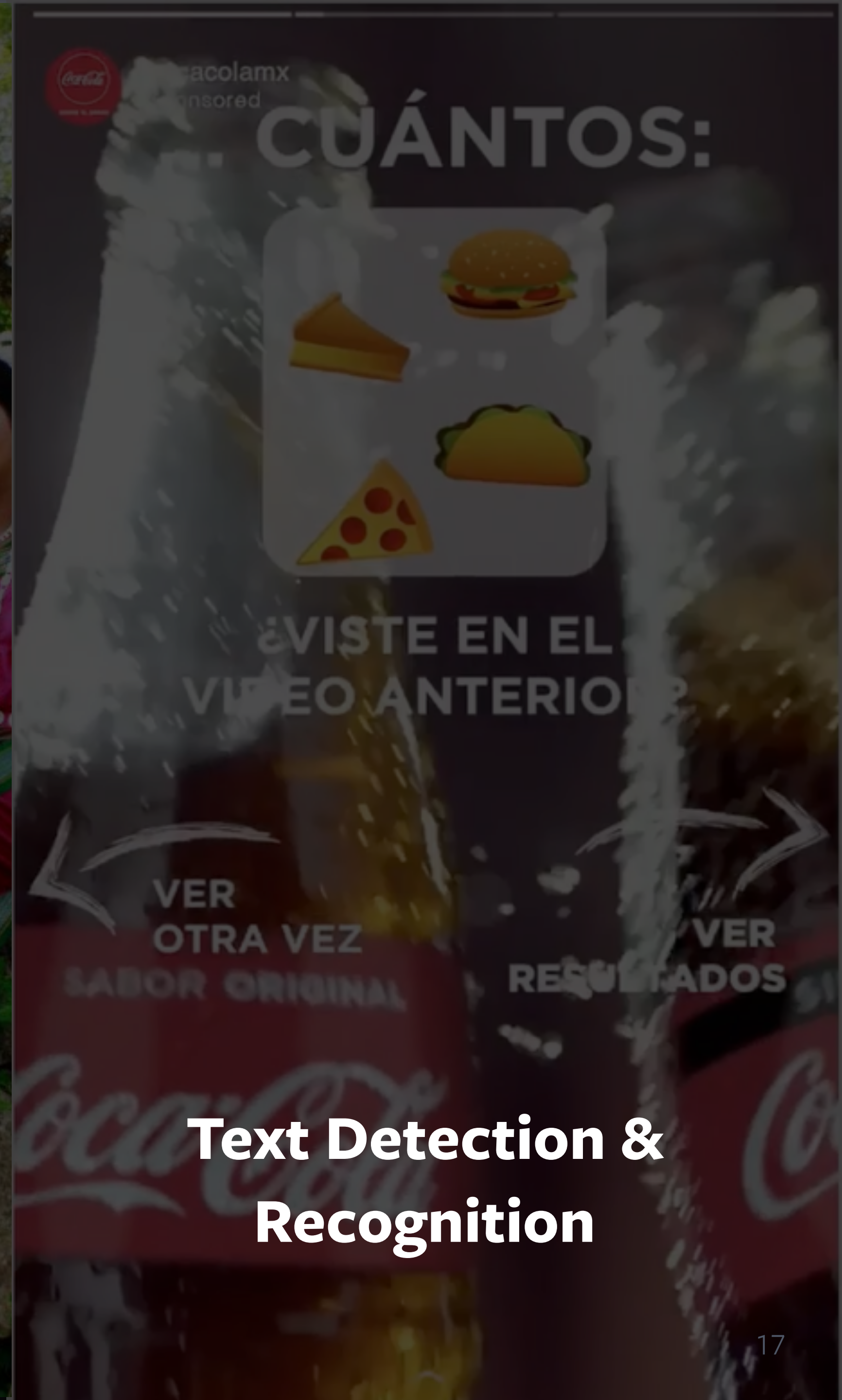




**Image & Video classification  
with thousands of concepts**

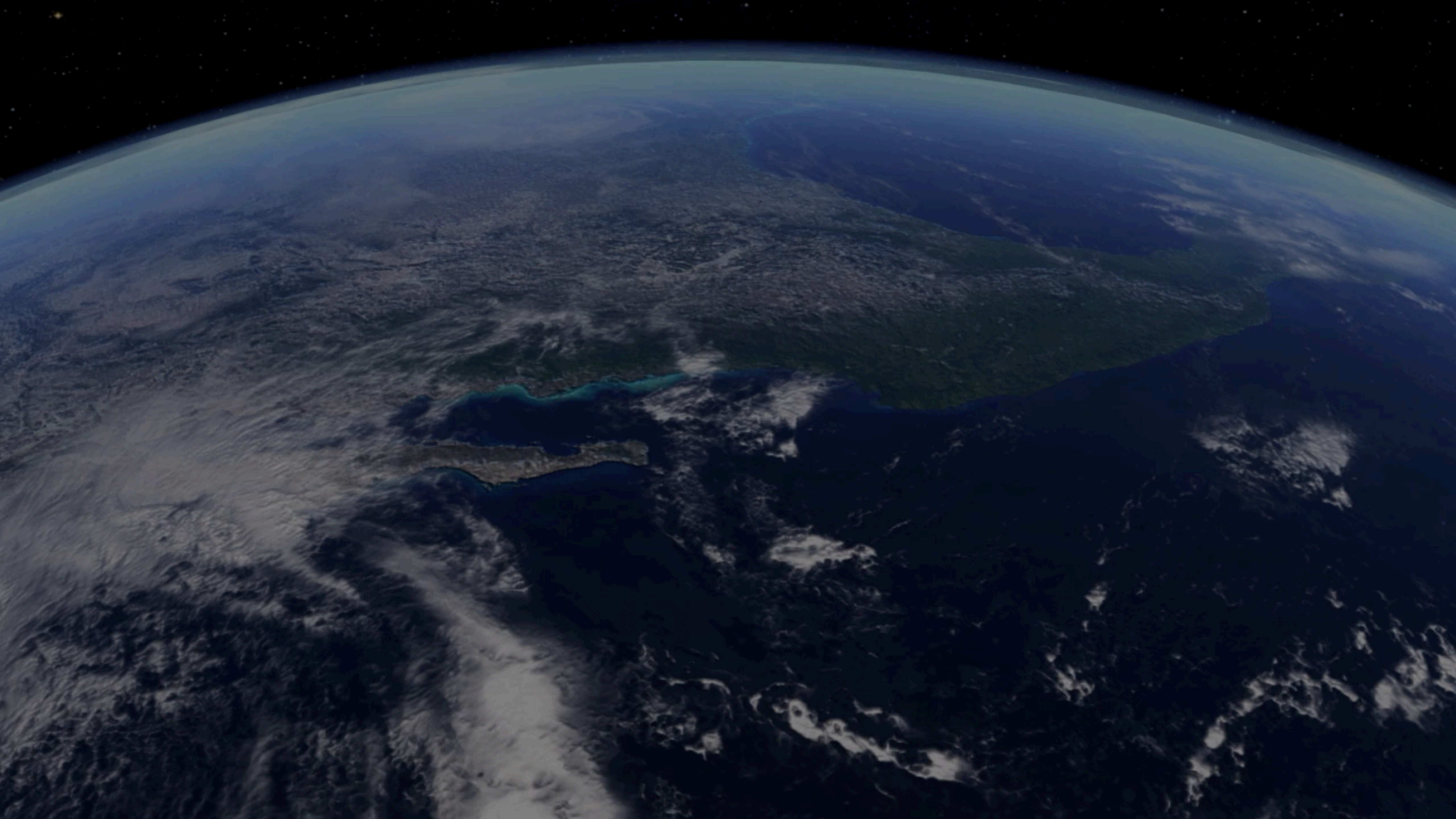


**Face & People of  
Interest recognition**



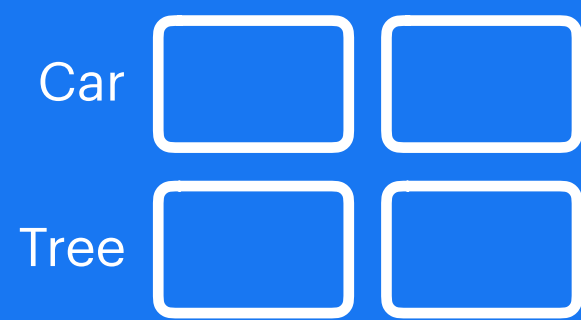
**Text Detection &  
Recognition**







# Levels of Supervision



Fully-Supervised



Weakly-Supervised



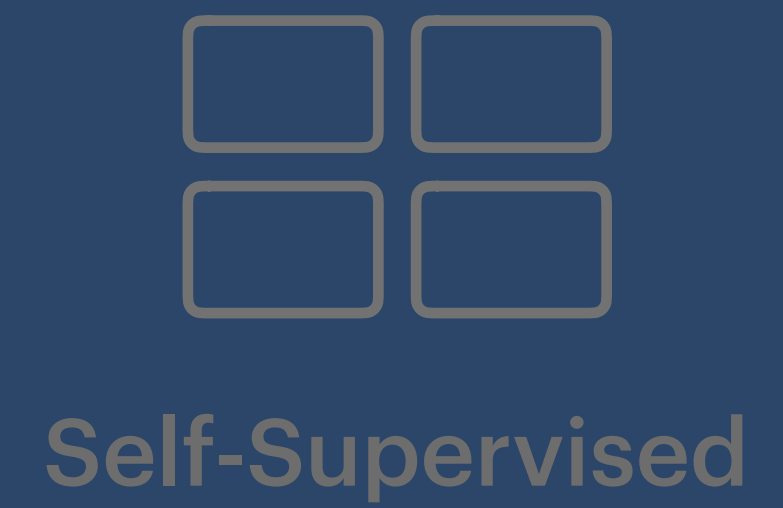
Semi-Supervised



Self-Supervised

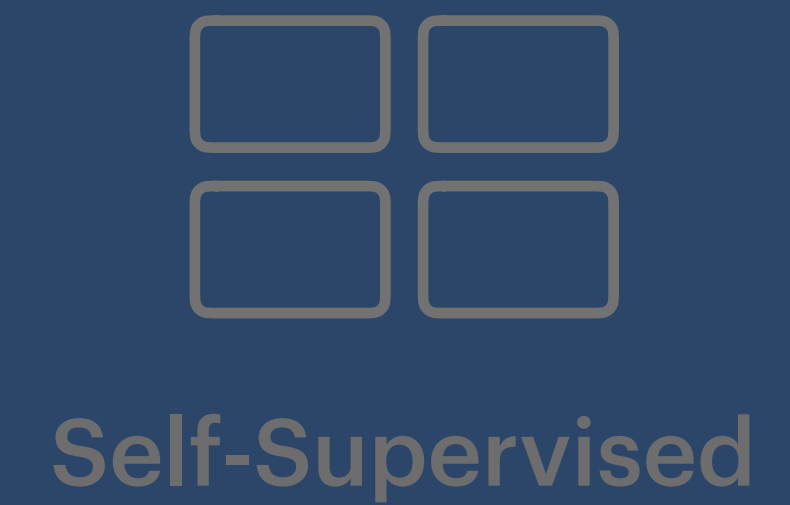


# Levels of Supervision



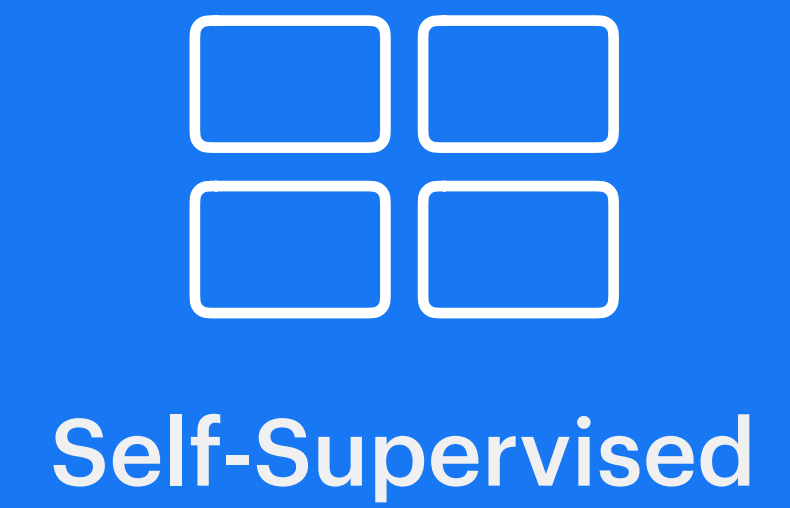


# Levels of Supervision





# Levels of Supervision

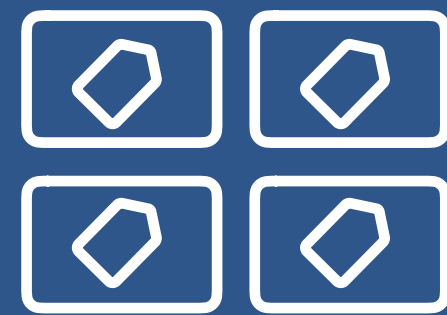




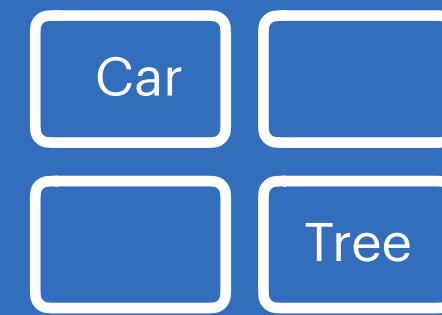
# Levels of Supervision



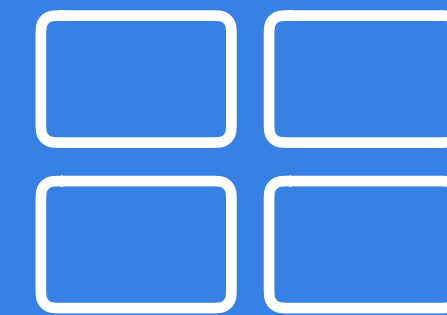
**Fully-Supervised**



**Weakly-Supervised**



**Semi-Supervised**



**Self-Supervised**







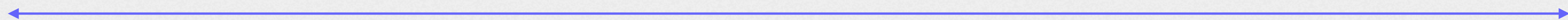
# Challenges of Training at Billion Scale

LEVELS OF SUPERVISION



Weekly supervised

Un-supervised



CAT, DOG, WOODEN FLOOR

???

A CUTE CAT COUPLE



# Challenges of Training at Billion Scale

NOISY DATA



Non-Visual Labels

#LOVE #CAT #DOG #HUSKY

Incorrect Labels

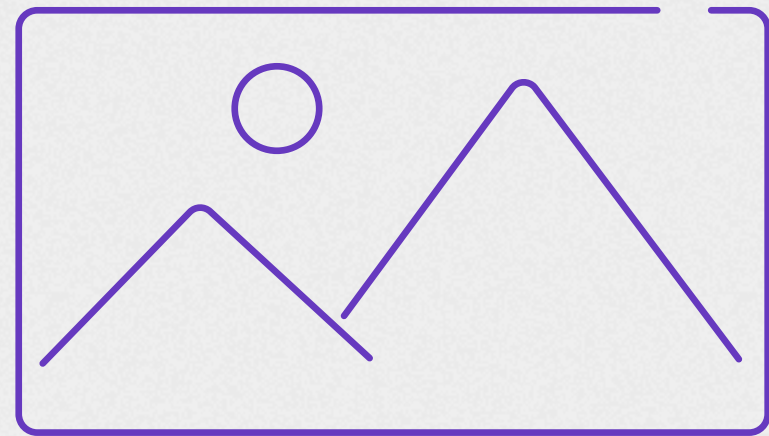
Missing Labels



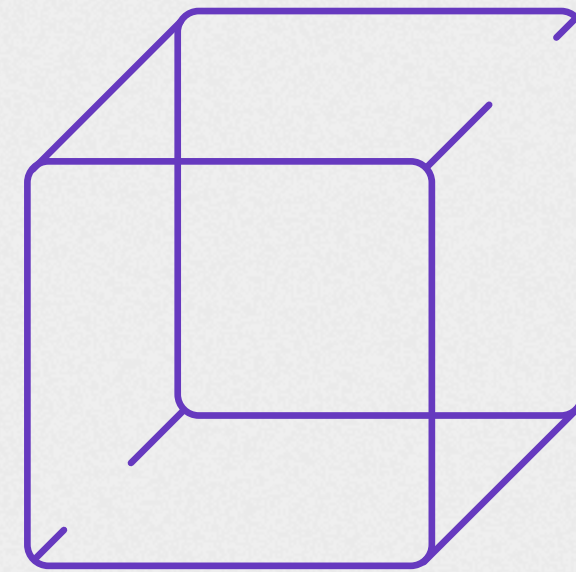




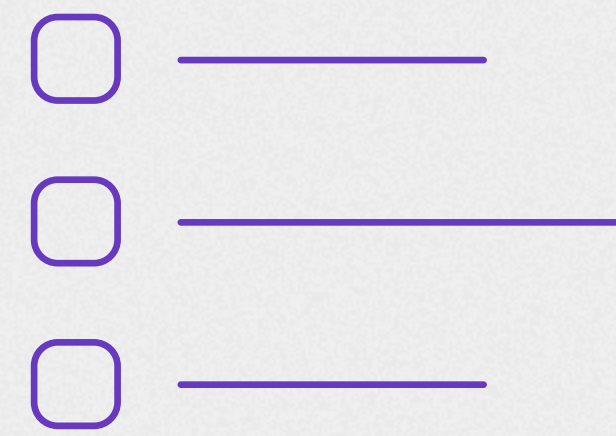
# Large Weakly Supervised Training



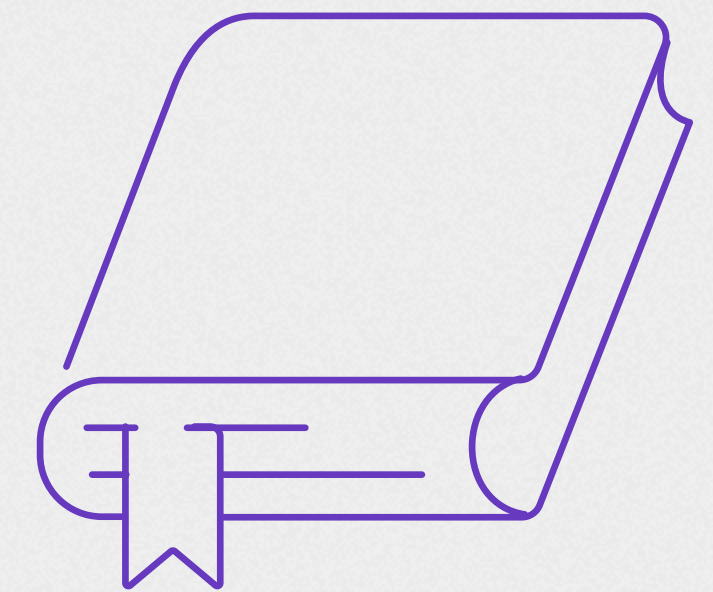
BILLIONS OF  
UNIQUE IMAGES



HUMUNGOUS  
MODELS



THOUSANDS OF  
LABELS



DISTRIBUTED  
TRAINING



# ImageNet in one hour

- ImageNet 1K has
  - 1.28 Million sample
  - 1000 categories
  - ResNet50 architecture
  - P100 GPUs
  - Caffe2

#machines	#gpus	Training time (mins)	Top-1 error
<b>1</b>	<b>8</b>	<b>1726.88</b>	<b>23.56</b>
2	16	905.22	23.35
4	32	464.23	23.48
8	64	231.42	23.39
16	128	117.76	23.58
<b>32</b>	<b>256</b>	<b>60.93</b>	<b>23.74</b>
36	288	54	23.76
40	320	48.84	24.08
44	352	44.36	24.23



# Billion Scale Training at FB

IMAGENET-1K: STATE OF THE ART RESULTS

**85.1%**

---

OUR 3.5B TRAINING  
RESNEXT101-32X32 MODEL

**83.1%**

---

PREVIOUS SOA



Before



Food



Food

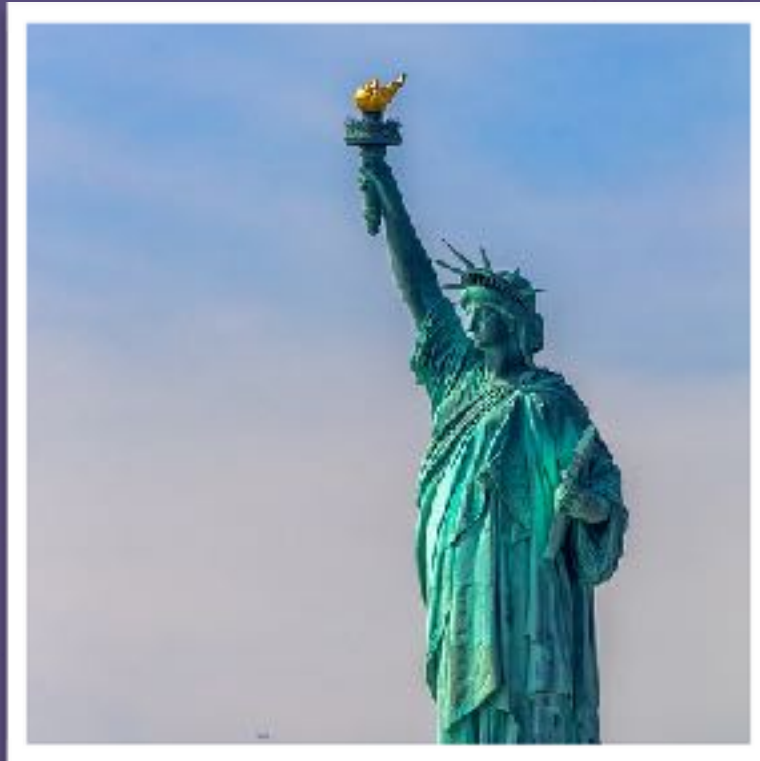
After



Cupcake



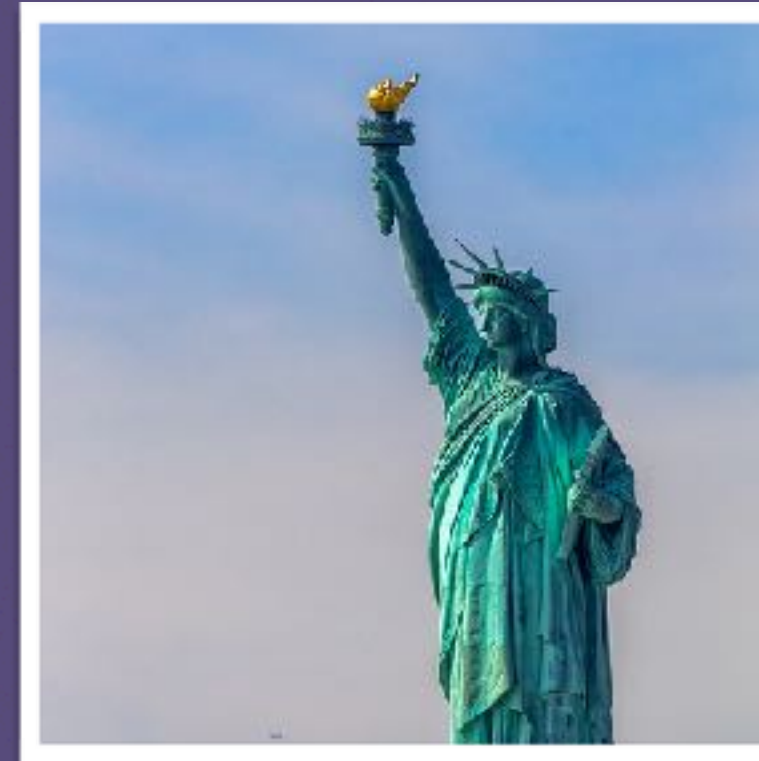
Apple pie



Landmark



???



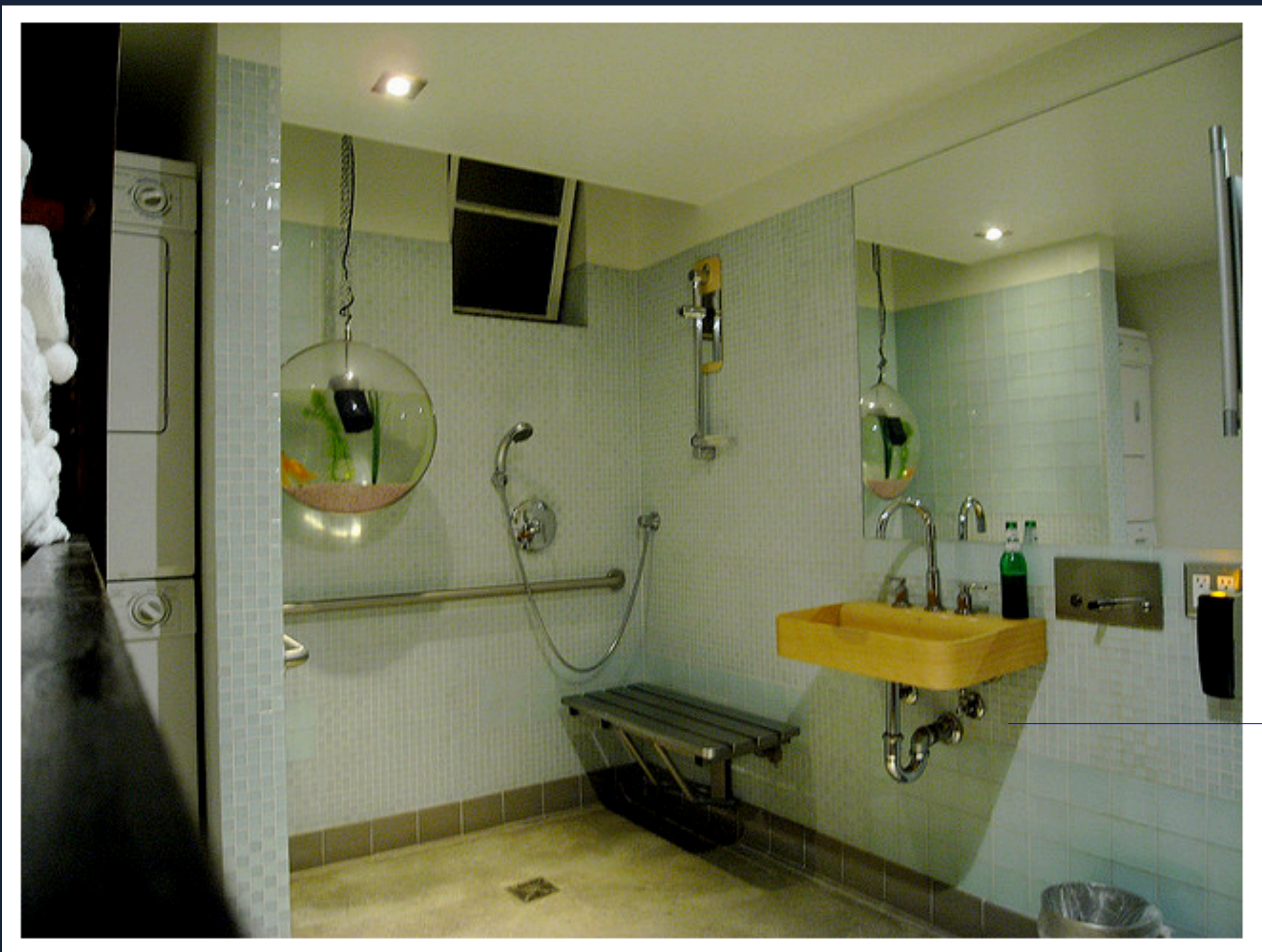
Statue of Liberty



Taj Mahal



WSL



IMAGENET

A simple, minimalist kitchen with very low lighting.

Black cat next to toy mouse on carpet.

WSL

A very modern bathroom with green glass tile work in the shower.

A cat sleeping next to a vehicle's tire on top of pavement.



# PyTorch Models Are Available

[https://pytorch.org/hub/facebookresearch\\_WSL-Images\\_resnext/](https://pytorch.org/hub/facebookresearch_WSL-Images_resnext/)

```
import torch
model = torch.hub.load('facebookresearch/WSL-Images', 'resnext101_32x8d_wsl')
# or
# model = torch.hub.load('facebookresearch/WSL-Images', 'resnext101_32x16d_wsl')
# or
# model = torch.hub.load('facebookresearch/WSL-Images', 'resnext101_32x32d_wsl')
# or
#model = torch.hub.load('facebookresearch/WSL-Images', 'resnext101_32x48d_wsl')
model.eval()
```



**How about Videos?**





---

**ENRICH OUR USER'S EXPERIENCE**



# Did We Miss a Great Moment?









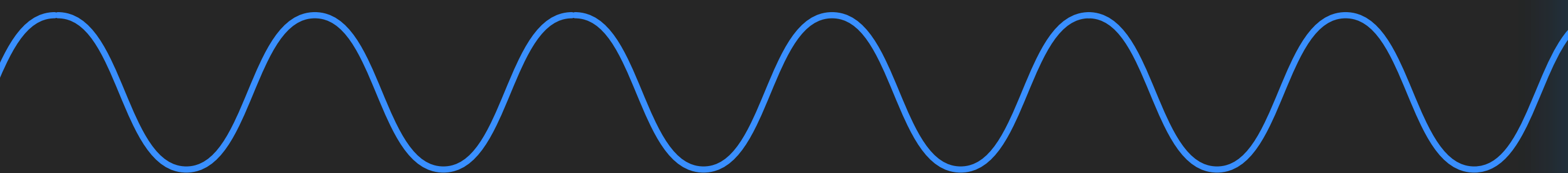
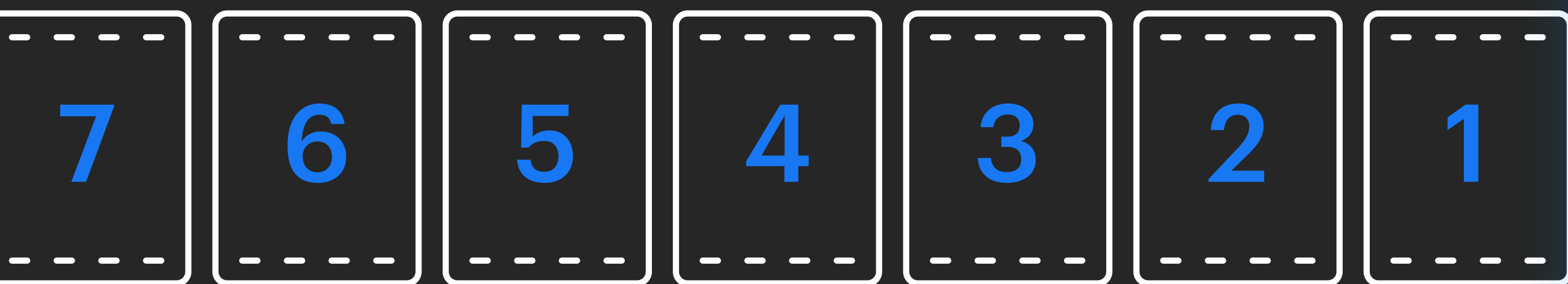
# Spatiotemporal Visual Modeling

SPACE AND TIME HAVE DIFFERENT STATISTICS





VIDEO

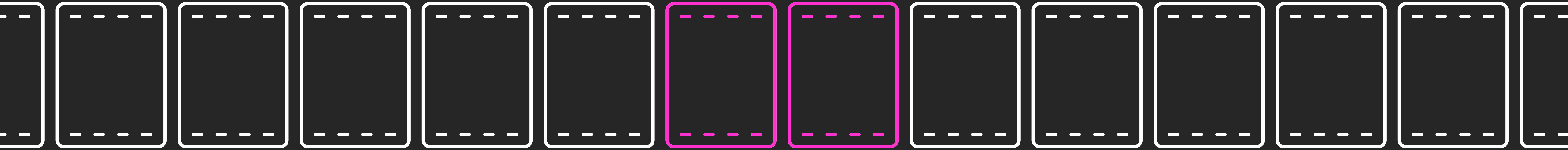


AUDIO





ACTION OF INTEREST

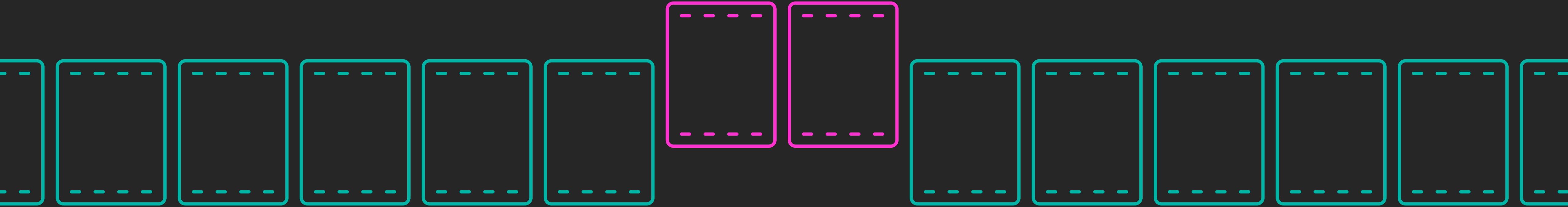






INFORMATION REDUNDANCY  
TEMPORAL NOISE

ACTION OF INTEREST























dog 0.99

door 0.8

wag 0.5

walk 0.8

stick 0.7

pickup 0.7

drop 0.8

tail 0.7





## Sampling Salient Clips

Spatiotemporal CNN

dog 0.99

door 0.8

wag 0.5

walk 0.8

stick 0.7

pickup 0.7

drop 0.8

tail 0.7







# Levels of Supervision

Car    
Tree  





**Fully-Supervised**  
~ Millions

**Weakly-Supervised**  
~ Billions

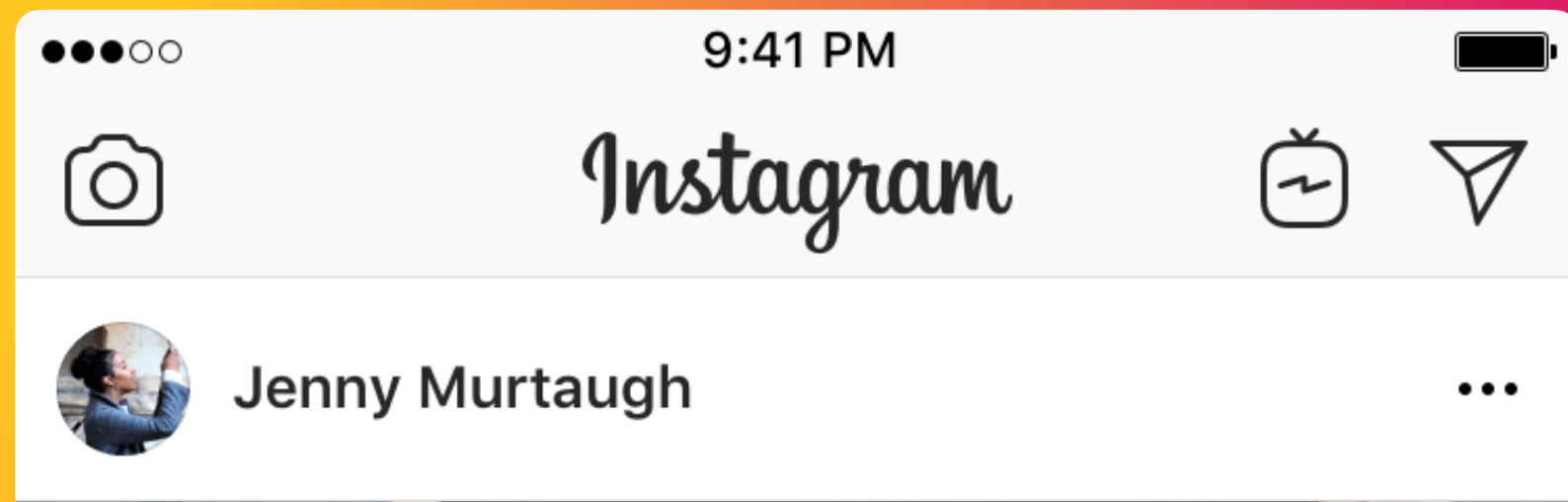
Car    
 Tree 

**Semi-Supervised**  
~ Trillions

**Self-Supervised**





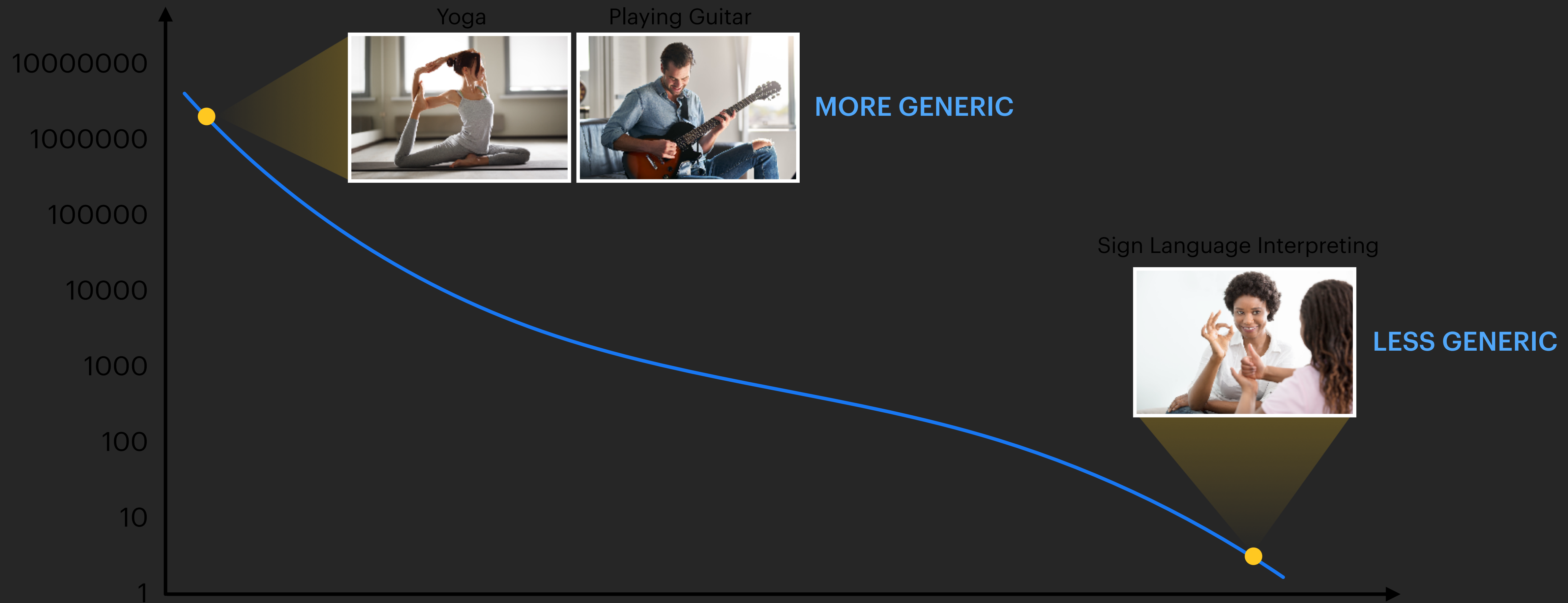
# Leveraging Rich Metadata

HASHTAGS



# Challenges With Training at Scale

## **SKEWED (LONG-TAIL) DISTRIBUTION**







65M

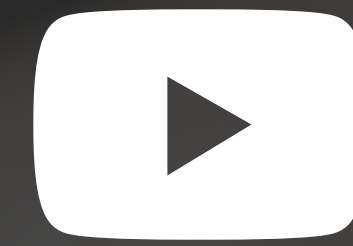
PUBLIC VIDEOS





65M

PUBLIC VIDEOS



6M

PUBLIC VIDEOS

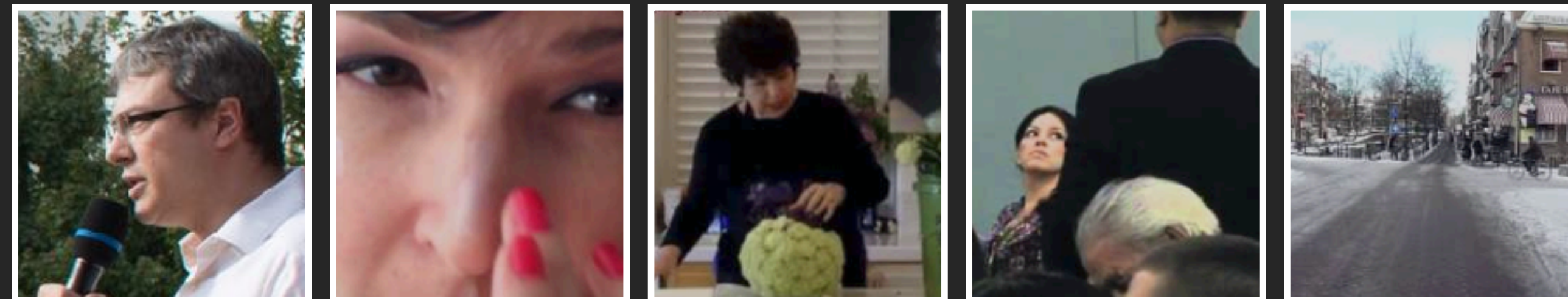


128M

PARAMETERS



# State of the Art Results



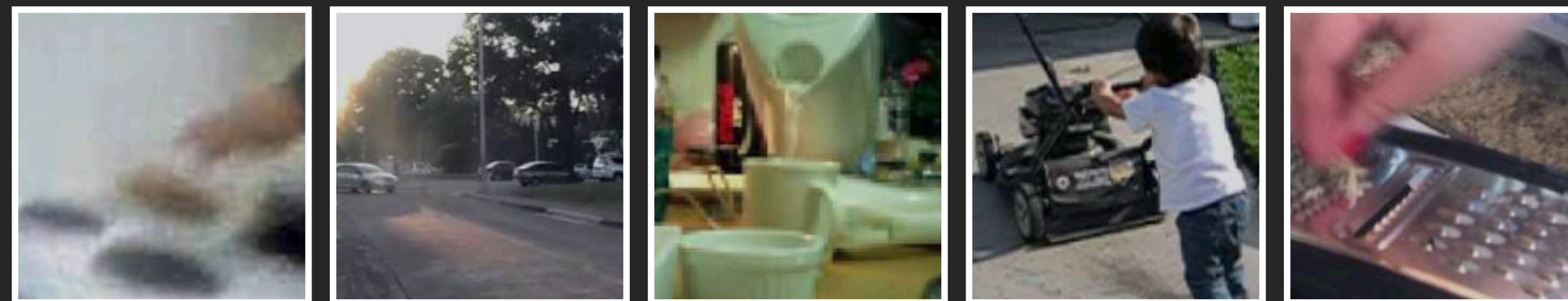
SPEAKING

APPLYING CREAM

ARRANGING FLOWERS

AUCTIONING

BIKING THROUGH SNOW



FLIPPING PANCAKE

JOGGING

MAKING TEA

MOWING

SCRAMBLING EGGS

Metric: Top-1 Accuracy

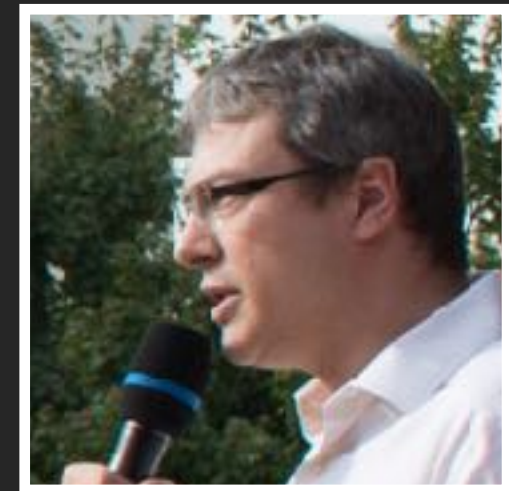
**77.7%**

PREVIOUS SOA

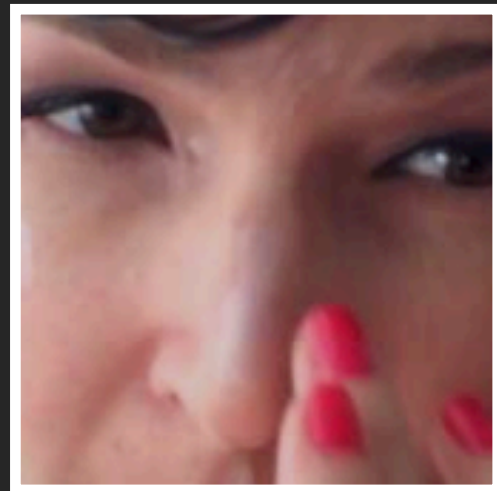
**KINETICS: 300K VIDEOS, 400 ACTIONS**



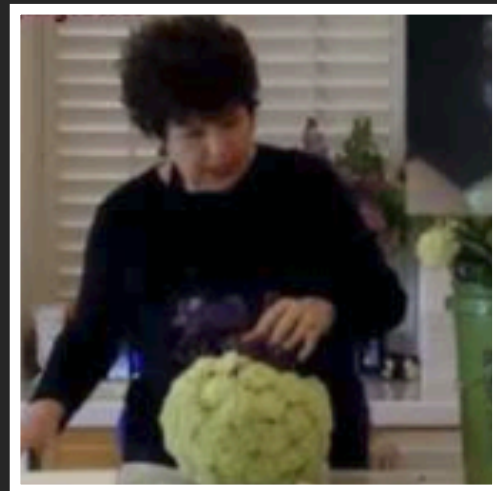
# State of the Art Results



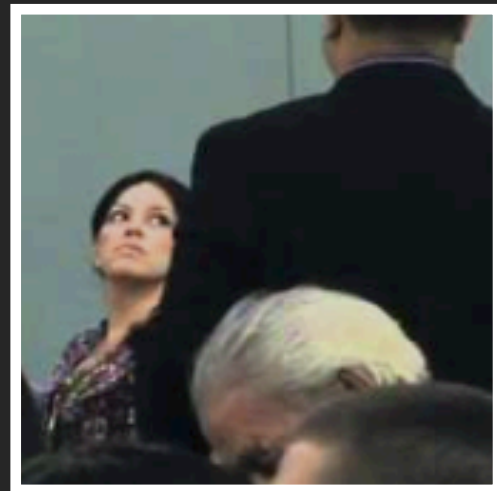
SPEAKING



APPLYING CREAM



ARRANGING FLOWERS



AUCTIONING



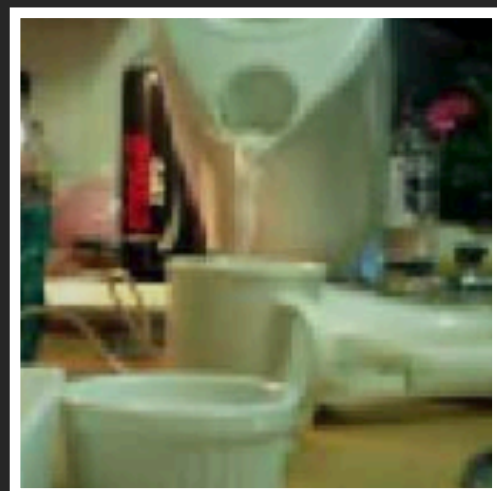
BIKING THROUGH SNOW



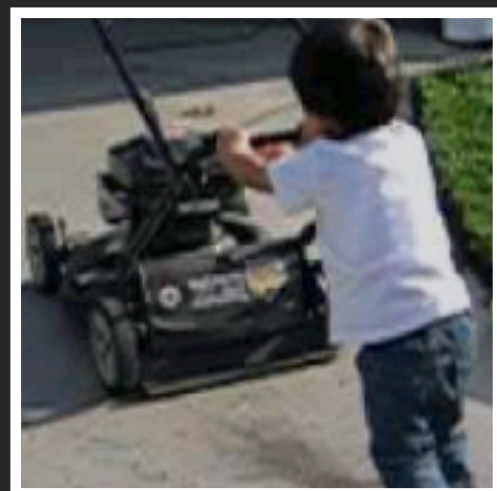
FLIPPING PANCAKE



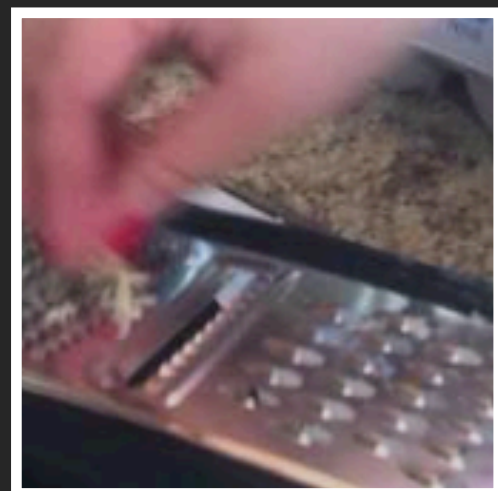
JOGGING



MAKING TEA



MOWING



SCRAMBLING EGGS

KINETICS: 300K VIDEOS, 400 ACTIONS

Metric: Top-1 Accuracy

77.7%

PREVIOUS SOA

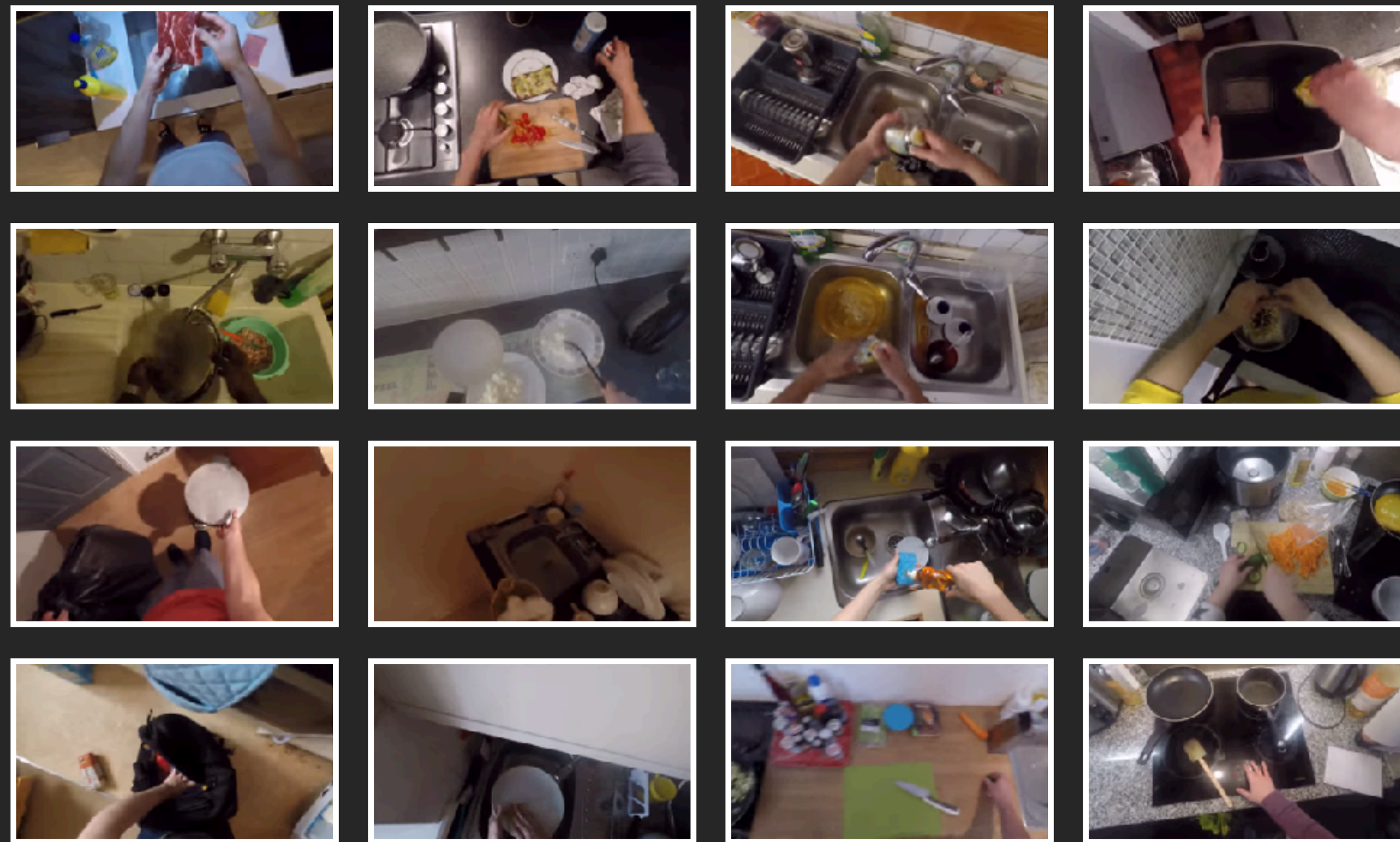
+5.1%

82.8%

OUR 65M TRAINING



# State of the Art Results



Metric: Top-1 Accuracy

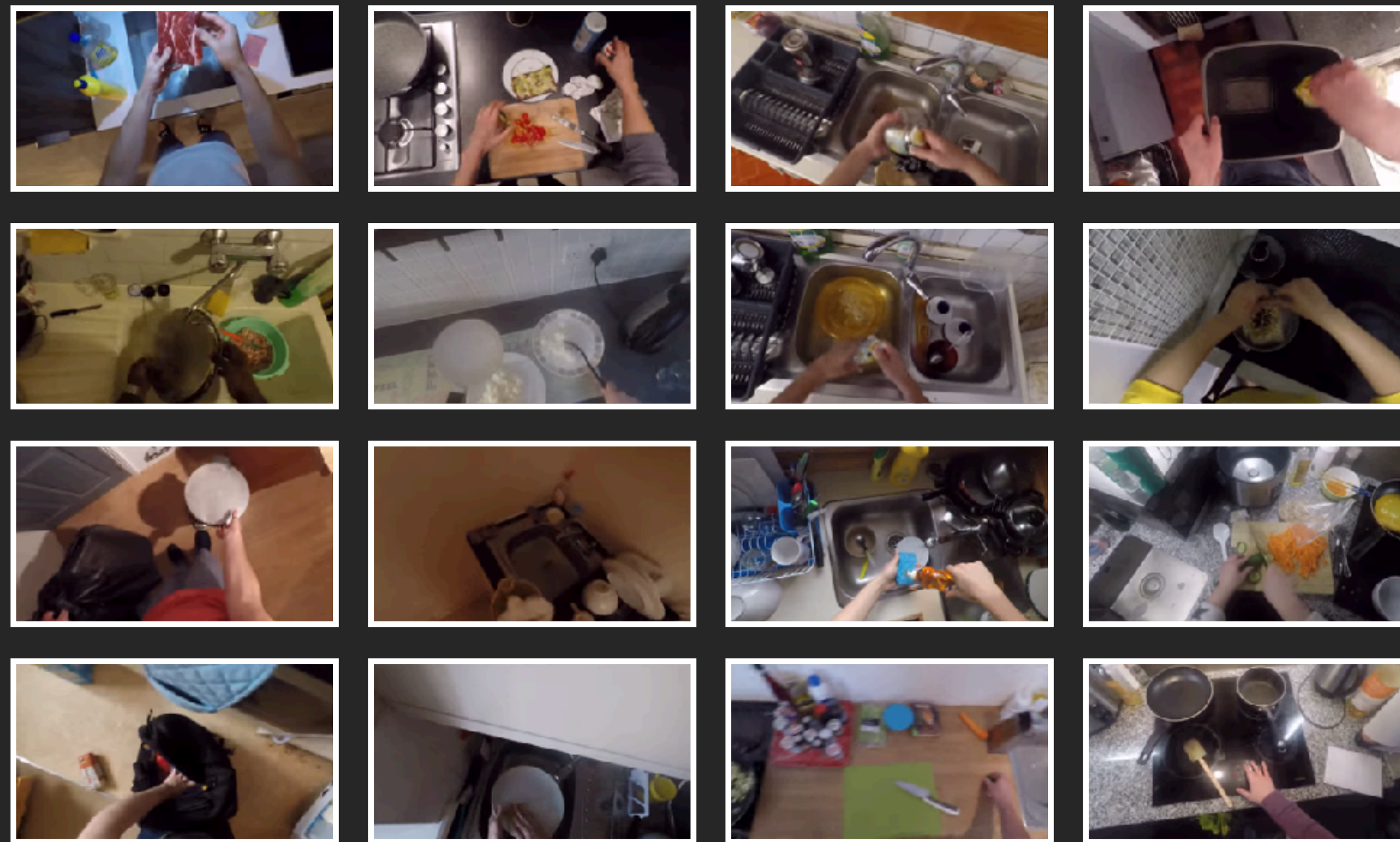
# 21.0%

PREVIOUS SOA

**EPIC-KITCHENS: 28K VIDEOS, 2337 ACTIONS**



# State of the Art Results



EPIC-KITCHENS: 28K VIDEOS, 2337 ACTIONS

Metric: Top-1 Accuracy

21.0%

PREVIOUS SOA

+4.6%

25.6%

OUR 65M TRAINING



# References for the video understanding efforts

- SCSampler: Sampling Salient Clips from Video for Efficient Action Recognition - <https://arxiv.org/abs/1904.04289>
- Video Classification with Channel-Separated Convolutional Networks - <https://arxiv.org/abs/1904.02811>
- Large-scale weakly-supervised pre-training for video action recognition - <https://arxiv.org/abs/1905.00561>




**Pushing State of the Art helps the world in significant ways**



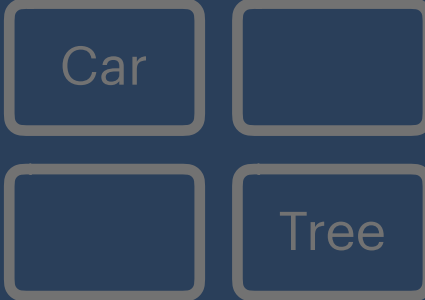
# Levels of Supervision

Car   
Tree 


**Fully-Supervised**  
~ Millions



**Weakly-Supervised**  
~ Billions



**Semi-Supervised**  
~ Trillions



**Self-Supervised**



## CrowdAI From Satellite Imagery to Disaster Insights

facebook

Jigar Doshi<sup>1</sup>, Saikat Basu<sup>2</sup>, Guan Pang<sup>2</sup>

CrowdAI<sup>1</sup>, Facebook<sup>2</sup>



### What's the Research?

- A framework for using convolutional neural networks (CNNs) on satellite imagery to identify the areas most severely affected by a disaster. This new method potentially produces **more accurate** information in **far less time** than current manual methods.
- The goal of this work is to allow rescue workers to quickly identify where aid is needed most, **without relying on manually annotated, disaster-specific data sets.**

### Disaster Impact Index (DII)

$$DII = \Delta Pred = \frac{\eta_{Pred_{before}=1 \& Pred_{after}=0}_{grid}}{\frac{1}{N_{grid}} \sum_{i=1}^{N_{grid}} \eta_{Pred_{before}=1}_{grid_i}}$$

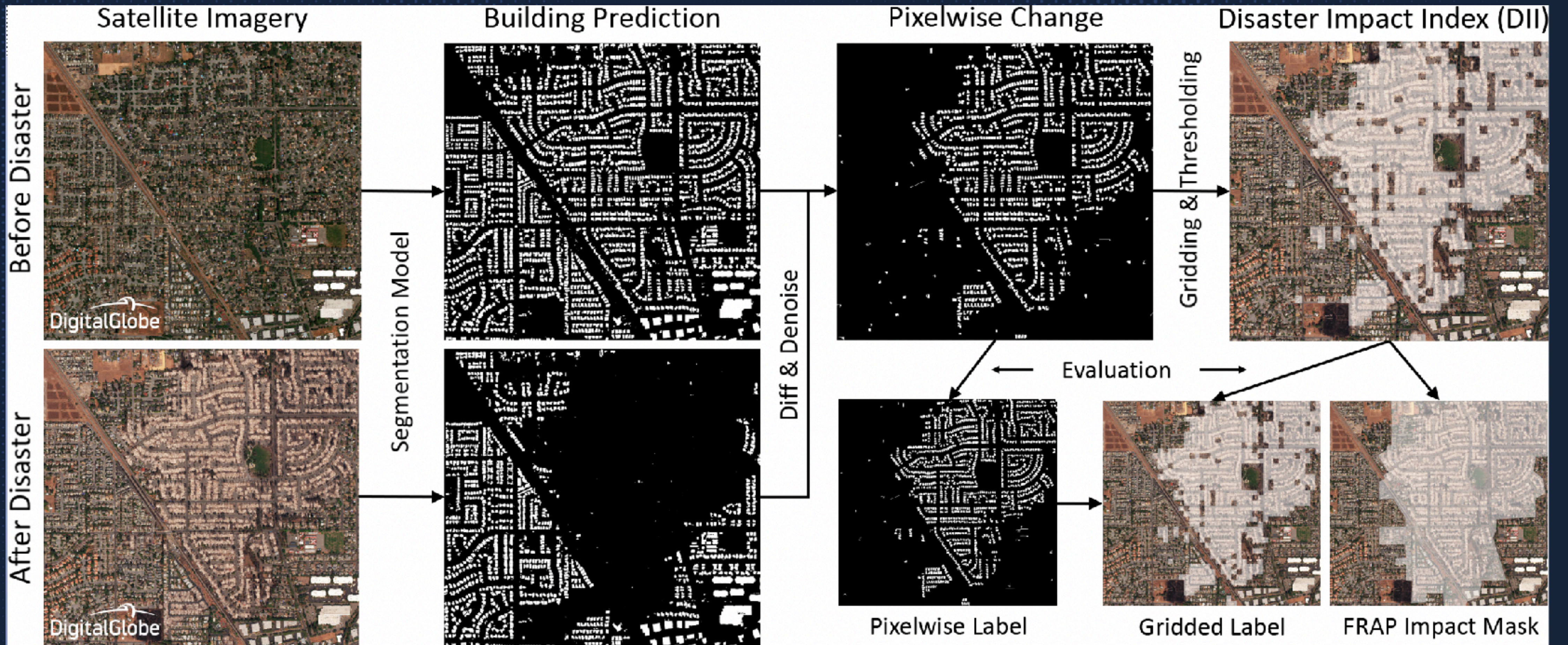
$\eta_{Pred_{before}=1 \& Pred_{after}=0}_{grid}$  → Pixels missing the feature post disaster

$\frac{1}{N_{grid}} \sum_{i=1}^{N_{grid}} \eta_{Pred_{before}=1}_{grid_i}$  → num of feature pixels predicted pre-disaster

$N_{grid}$  → total number of grids in this case  $256 \times 256$



# How it works!





# Helping in real world scenarios

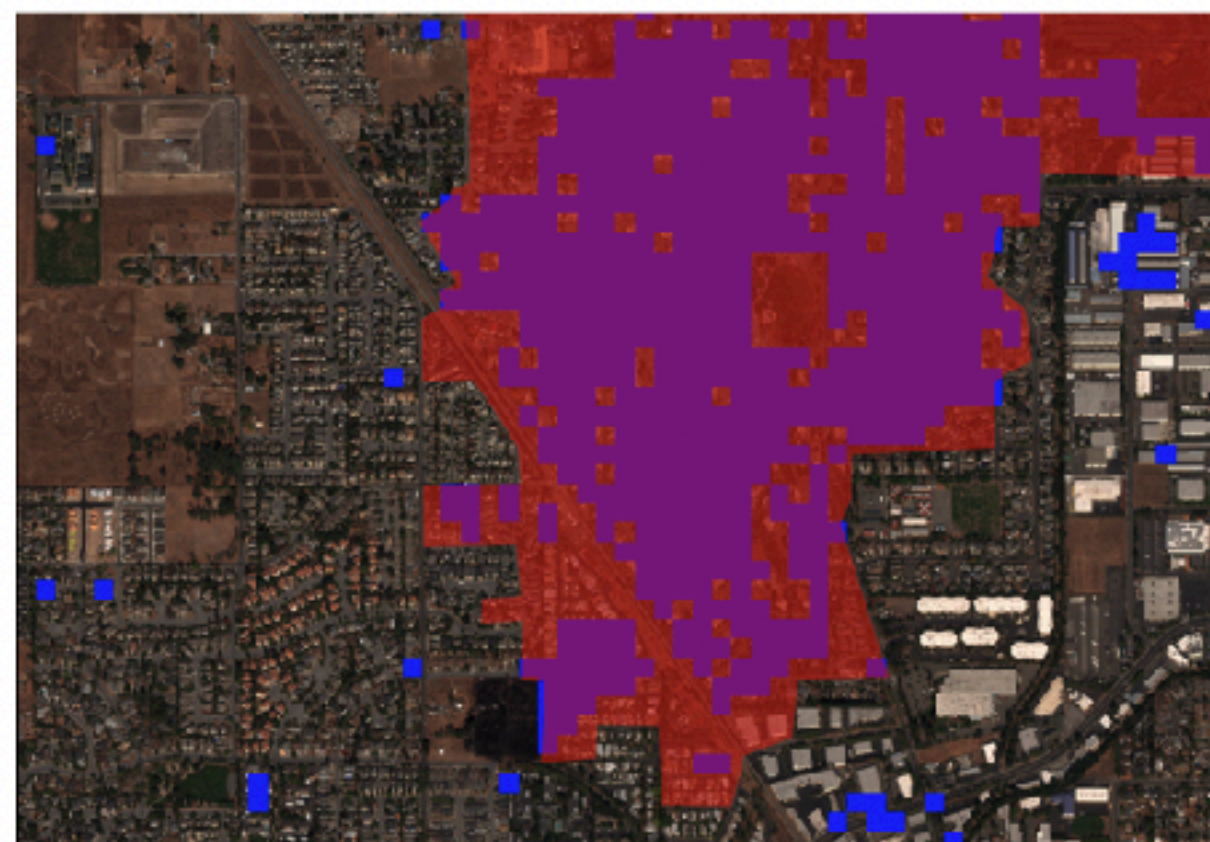
## Santa Rosa Fire - Buildings



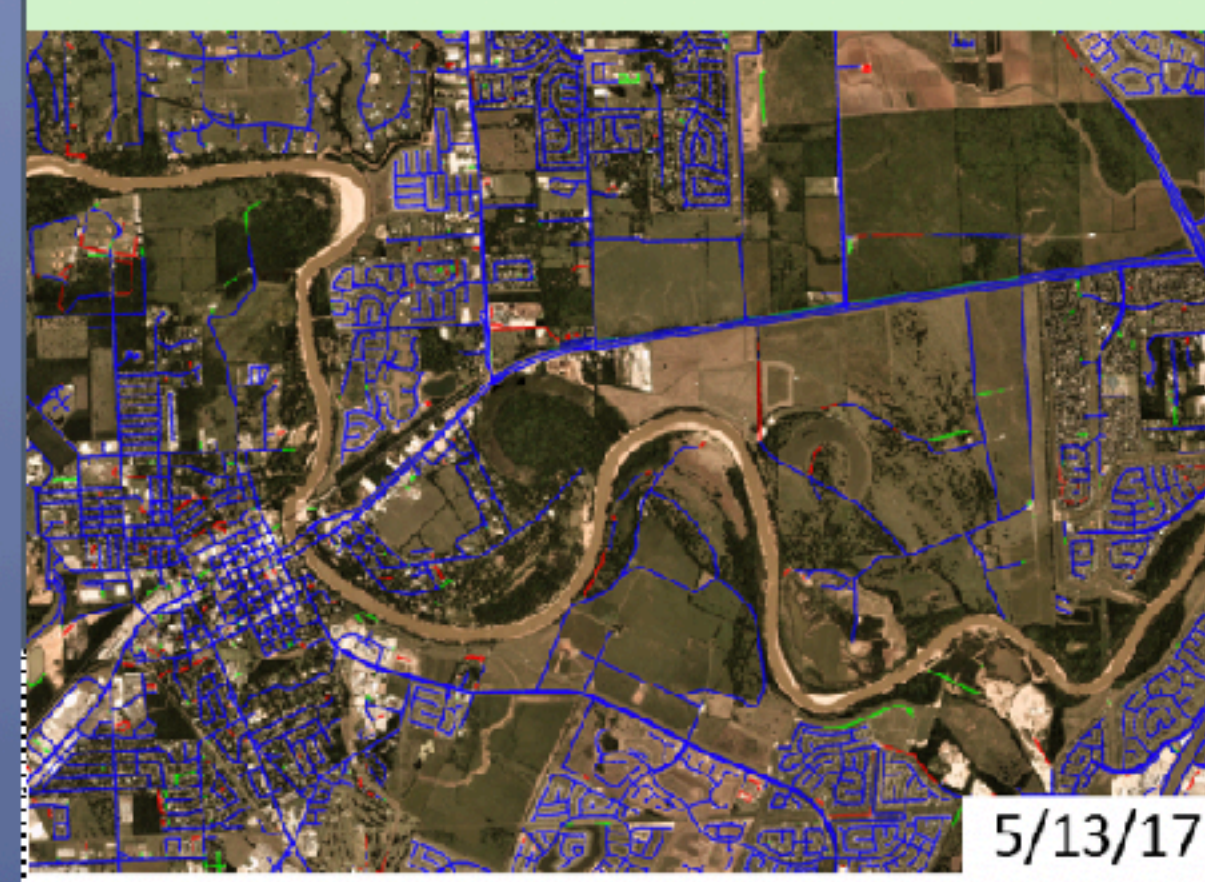
(a)



(b)

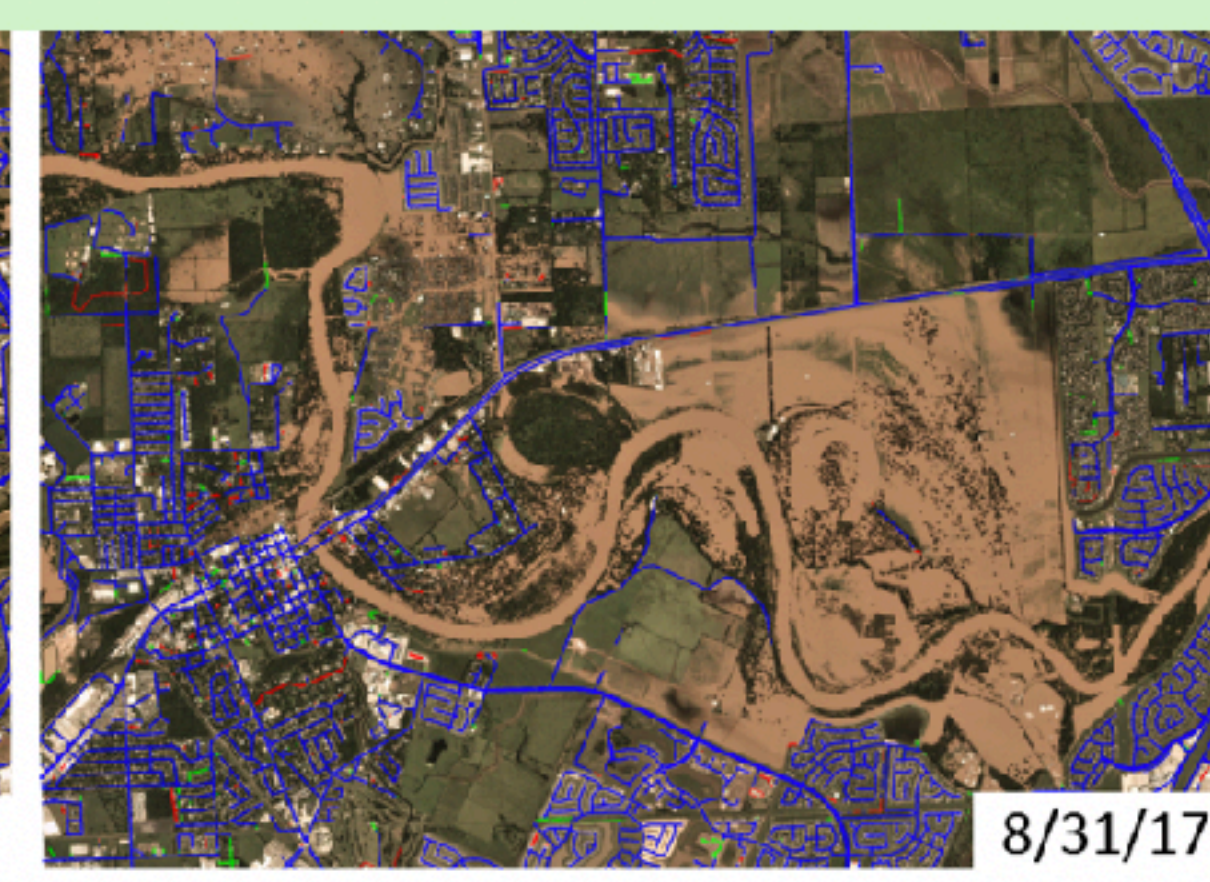


## Hurricane Harvey Flooding - Roads



5/13/17

(a)



8/31/17

(b)





# 1st CVPR Workshop on COMPUTER VISION FOR GLOBAL CHALLENGES

16/17 June 2019  
Long Beach, CA



Computer Vision  
for Global  
Challenges  
@cv4gc

Do you have an idea for a computer vision task that would impact the lives of many? Have you identified the limitations of a vision technique because of the geographical bias of the data you are using it on? Is there an application of computer vision that would be helpful to your community? Are you looking for potential vision expert partners or feedback on your idea?

Apply for the Call of Challenges, and come and experience the premier computer vision conference, and participate in an active discussion with the top vision researchers!

More info at: <http://www.cv4gc.org/#challenges>



# Population Density Estimation





# Naivasha in Kenya



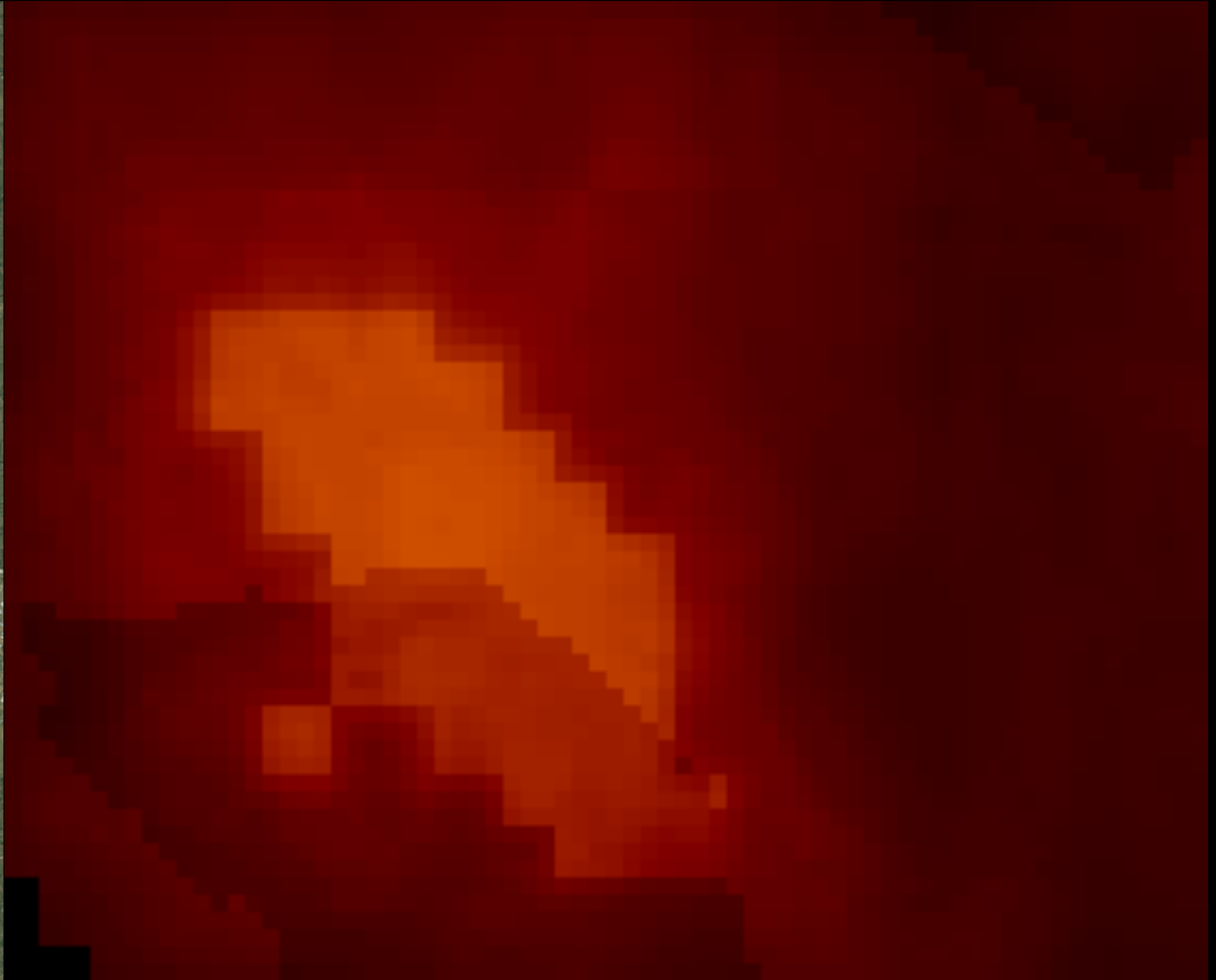


# GPWv4 from CIESIN at Columbia University



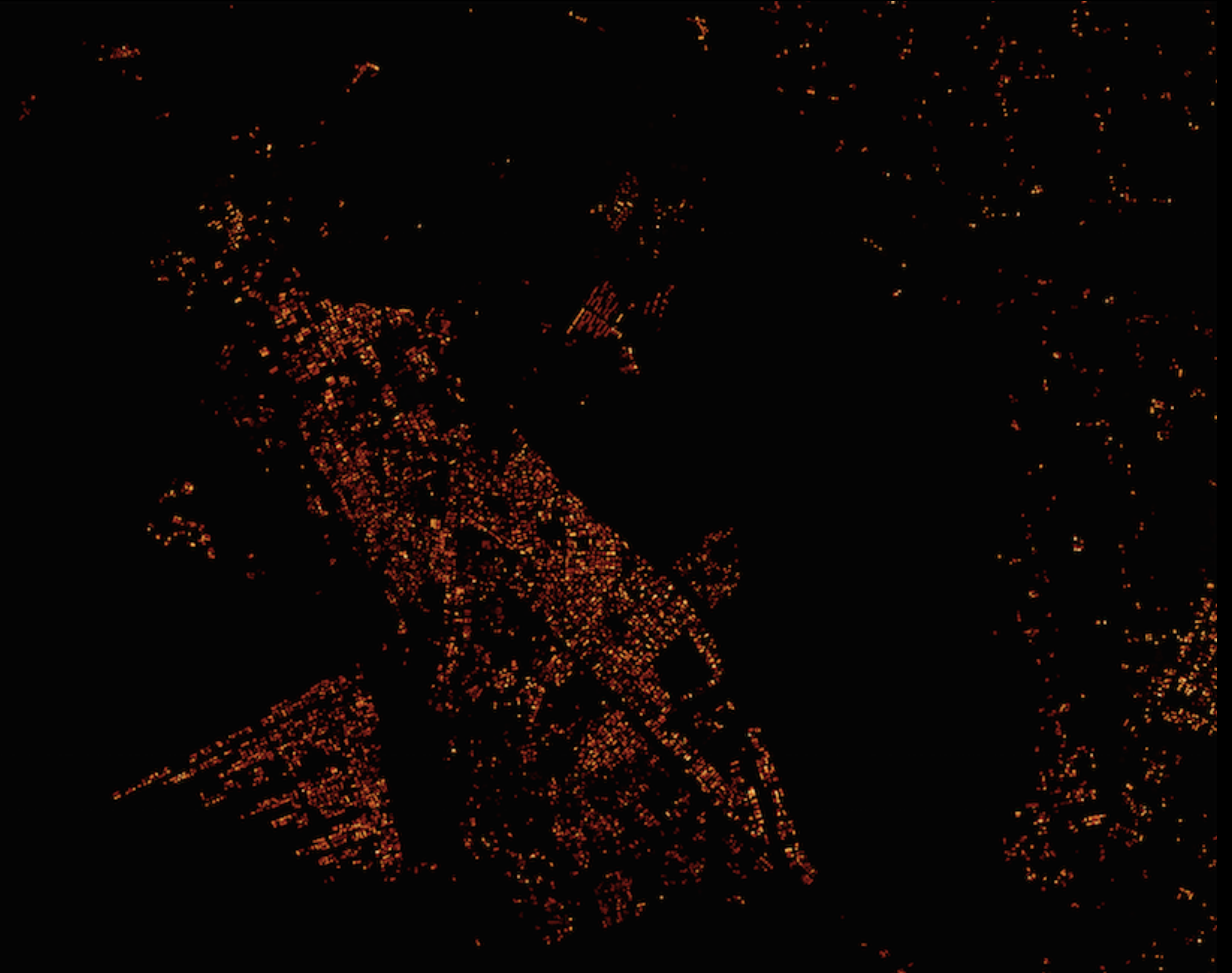


# WorldPop





# Facebook





What will you learn today?

How do you design a image and video recognition system for billion scale?

Can you remove the requirement of annotation to learn best representations?

Can we understand video faster than understanding individual frames?

How does pushing state of the art in CV make a meaningful difference to everyone in the world?



**Thank You!**

