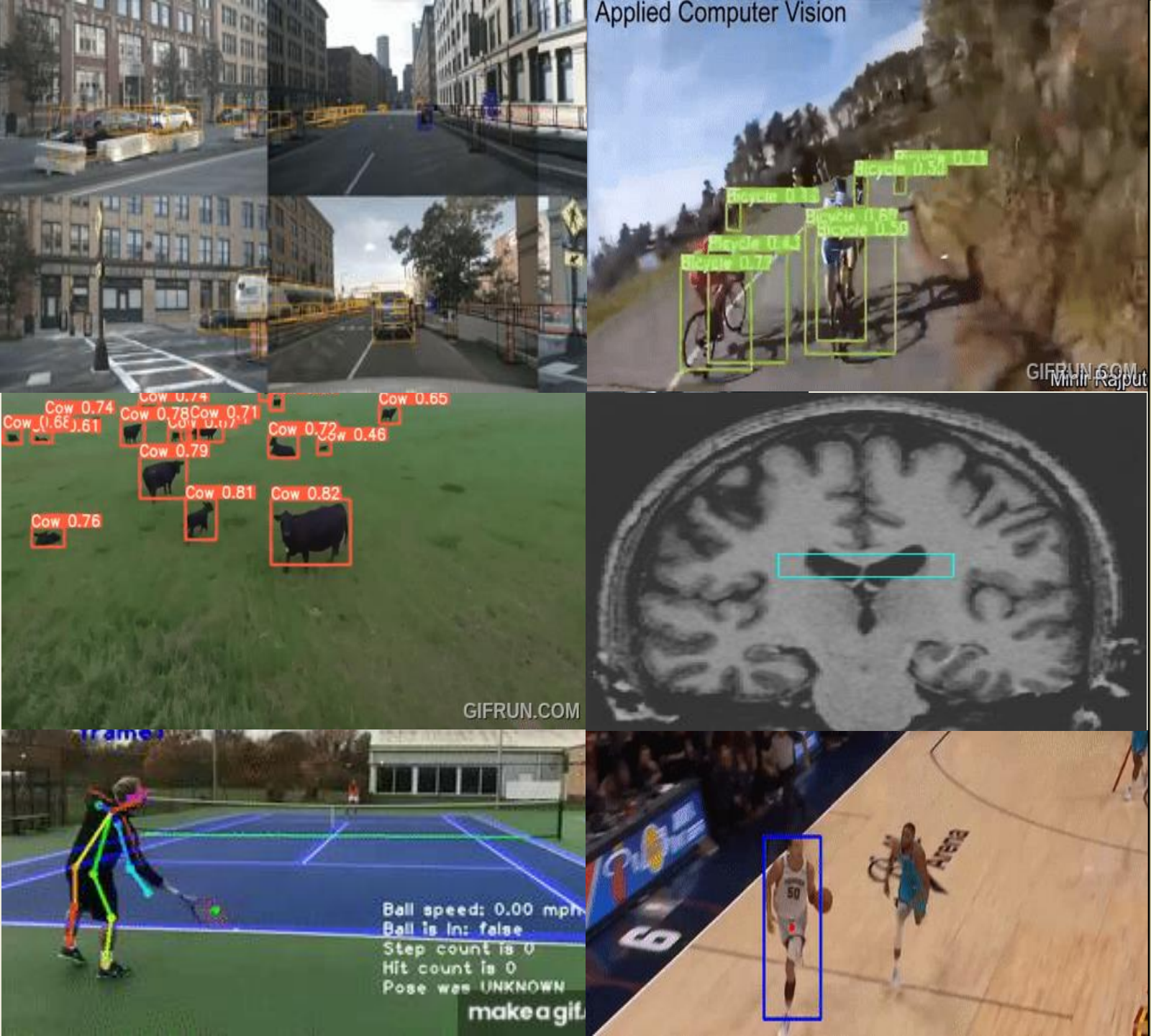


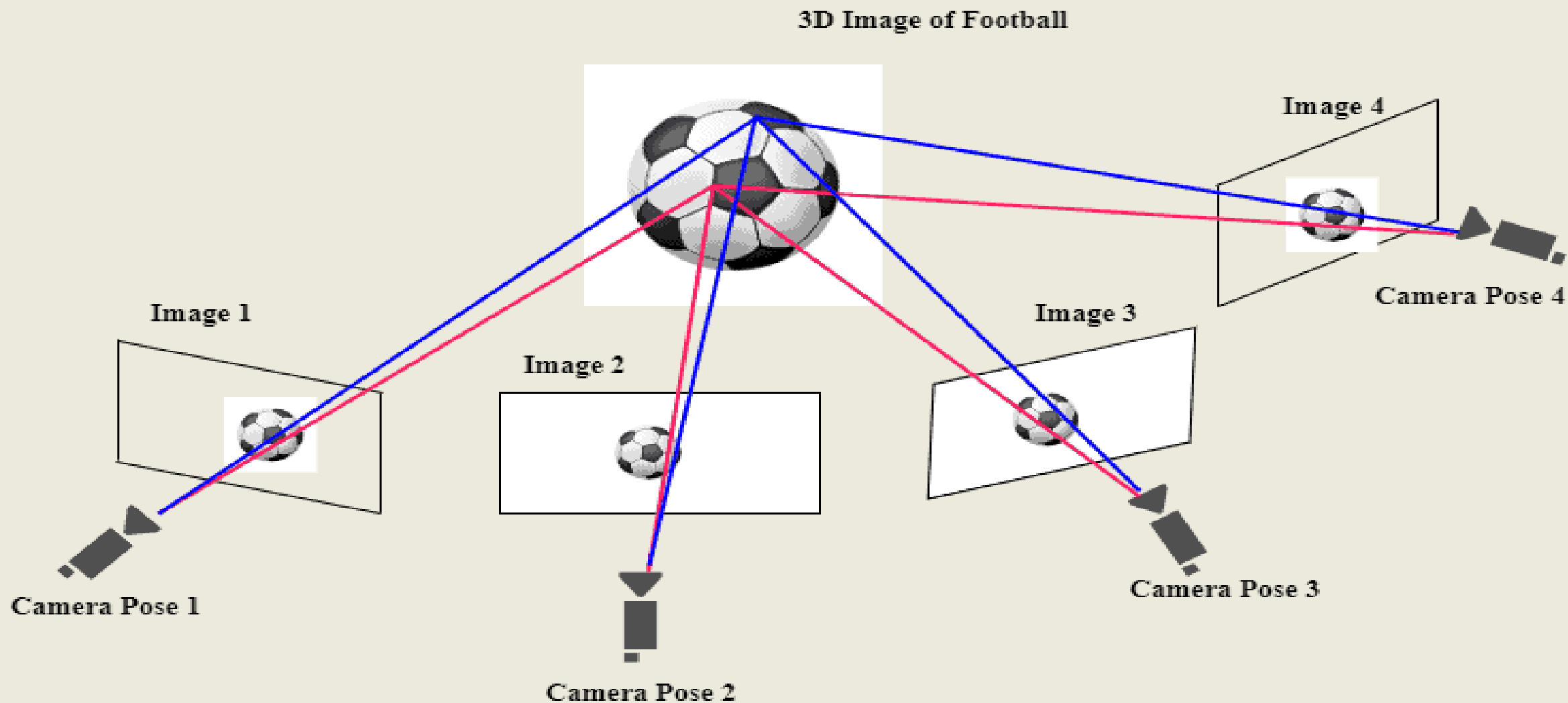
Multi-View Obscured Object Detection (MOOD)

Robin Arun, Katie Hucker, Afraa Noureen, Jiayi Zhou

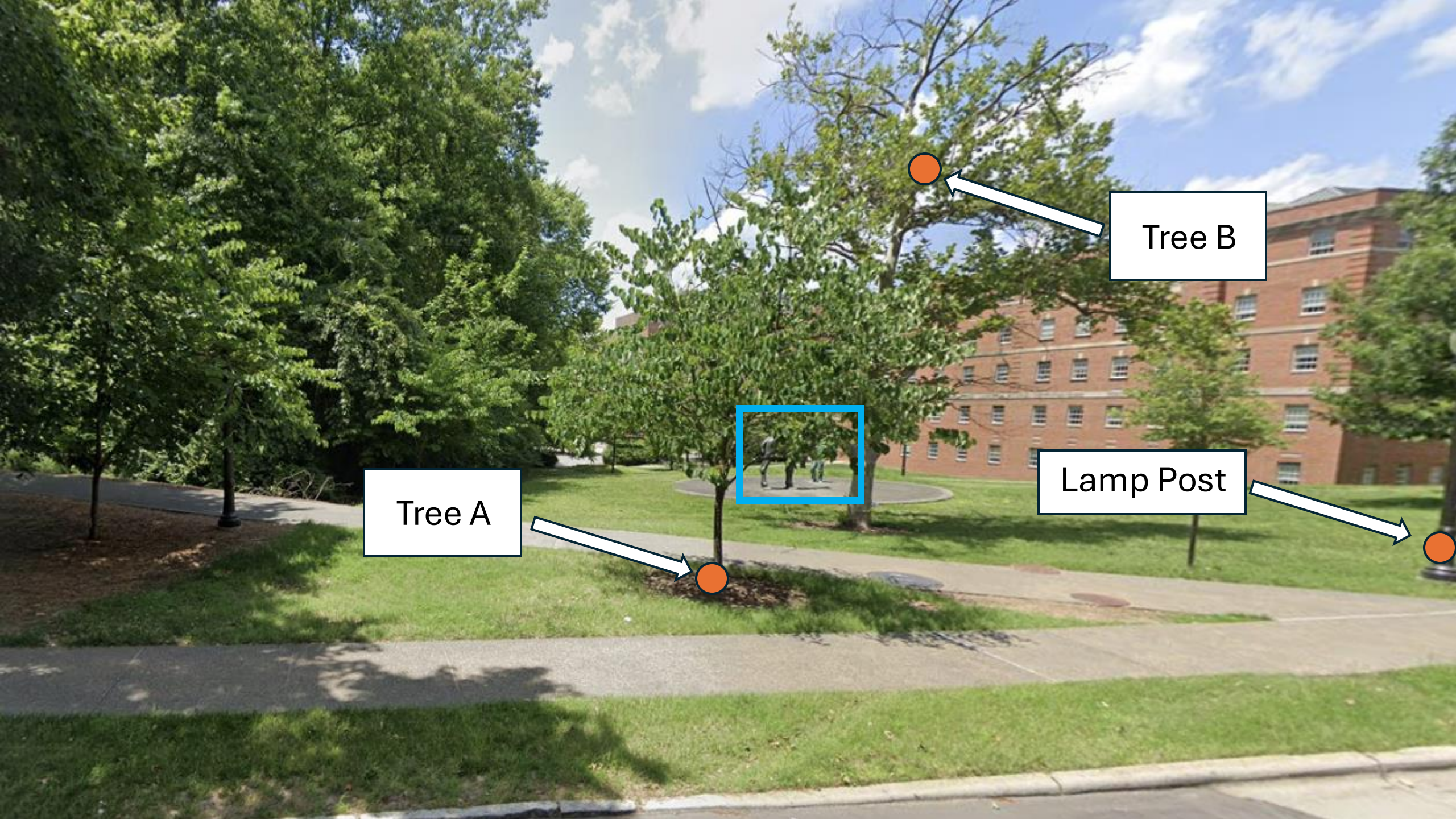


Introduction:

- Life in a 3D world
- Pedestrians, cars, animals.
- Movement, depth, visibility, light
- Pictures and video in a 2D plane.



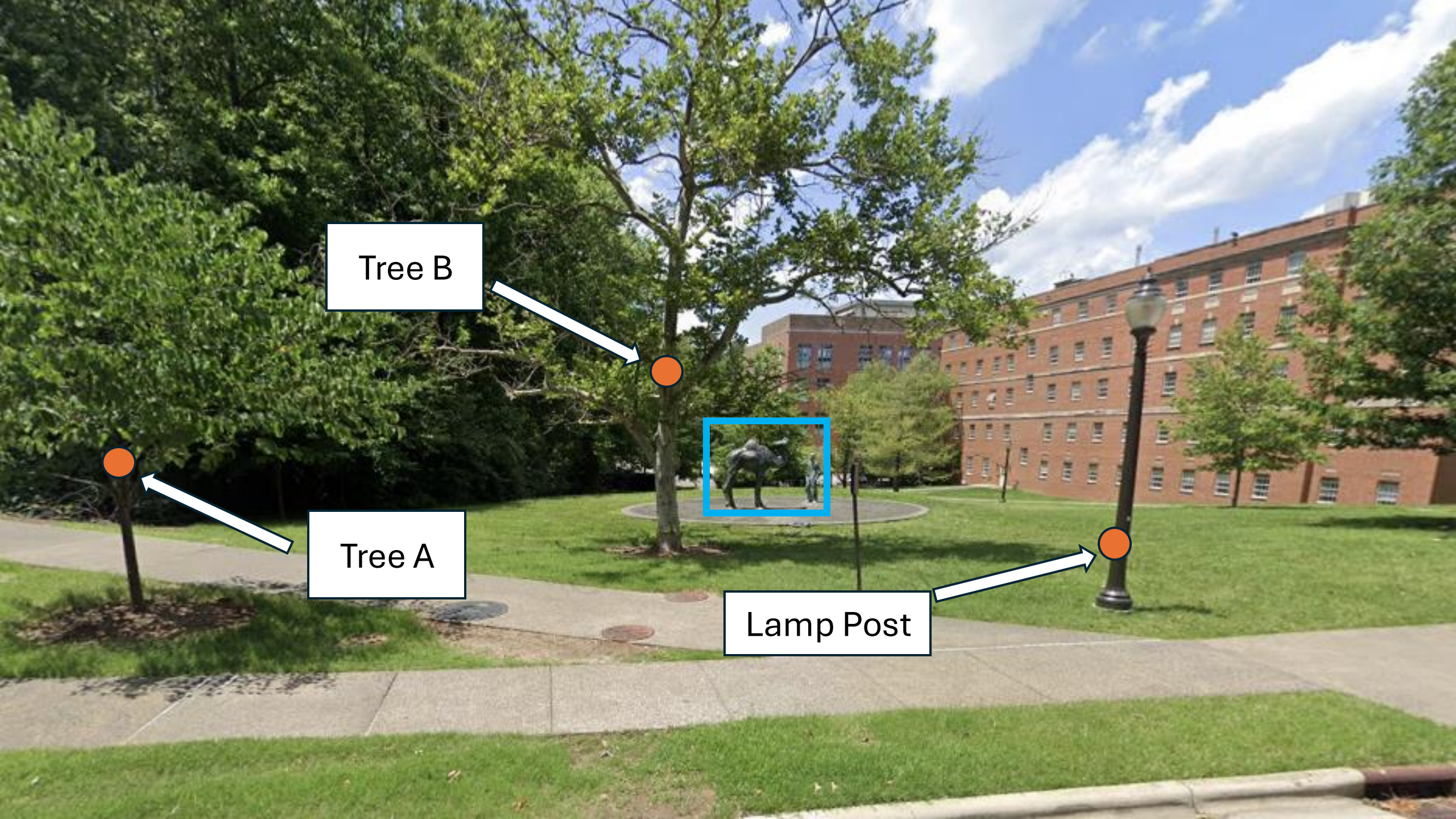
Multiple images or *views* can create a 3D world



Tree A

Tree B

Lamp Post



Tree B



Tree A



Lamp Post



Definition

- Multi-view: multiple images of the **same** scene
- Obscured object: Targets within a scene are difficult to identify due to the angle of view and scene debris like trees, people, and buildings.
- Detection: Can we identify if the object is there, and what the object is? How confident are we?



Our Problem

The image is a composite of two aerial photographs of a city, separated by a horizontal green banner. The top photograph shows a large, dark, rectangular building on the left and a bridge structure on the right. The bottom photograph shows a large number of small, black, geometric shapes, resembling wireframe cubes or pyramids, scattered across the lower half of the image. The banner in the center contains the text "Our Problem" in a large, black, sans-serif font.

What is the problem?

- We **sacrifice** scene contextualization and poor object visibility from single 2D images and videos, *unless we have one good look.*
- A model which **understands the 3D world** using co-registered images can provide stronger detection capabilities, *using **all the looks.***

Why does this model matter?

- Hidden Objects are a **risk**.
- **High-Stakes** impact
- Traditional methods are computationally expensive.
- A 3D interpretation of the world without a 3D reconstruction

Goal: Develop a model which can...

1. **input** multiple pictures from the same scene
2. **use** how those pictures relate to one another within the scene
3. **detect** and classify the obscured object of interest *correctly*
4. **interact** in a useable tool

How much can partially
obscured object
detection improve using
multi-views?

1.

How do different levels of
obscurement with
multiple views affect
detection accuracy?

2.

What number of scene
views do we need for
high detection
performance?

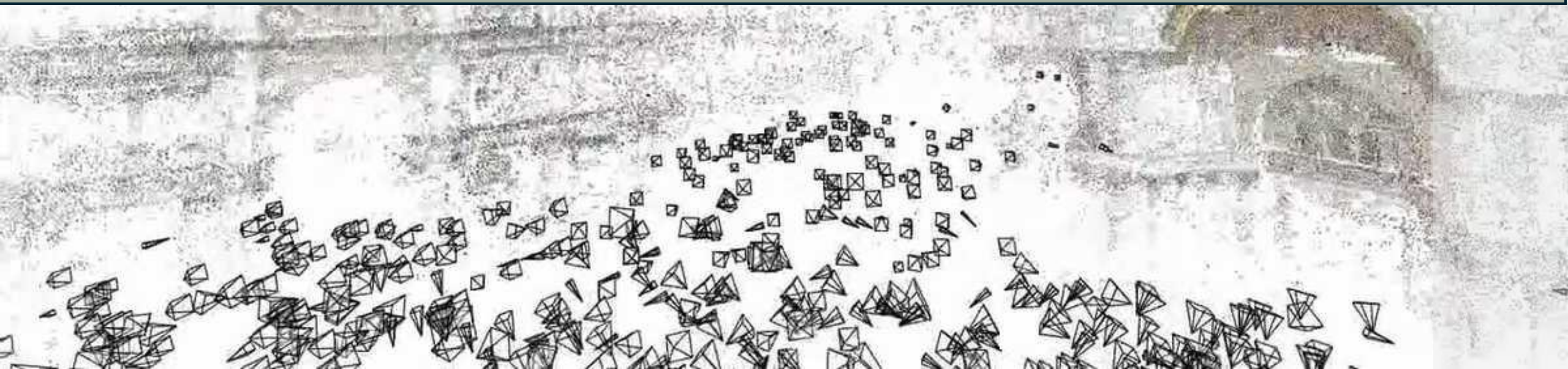
3.

What features of
our objects
contribute to
detected boxes?

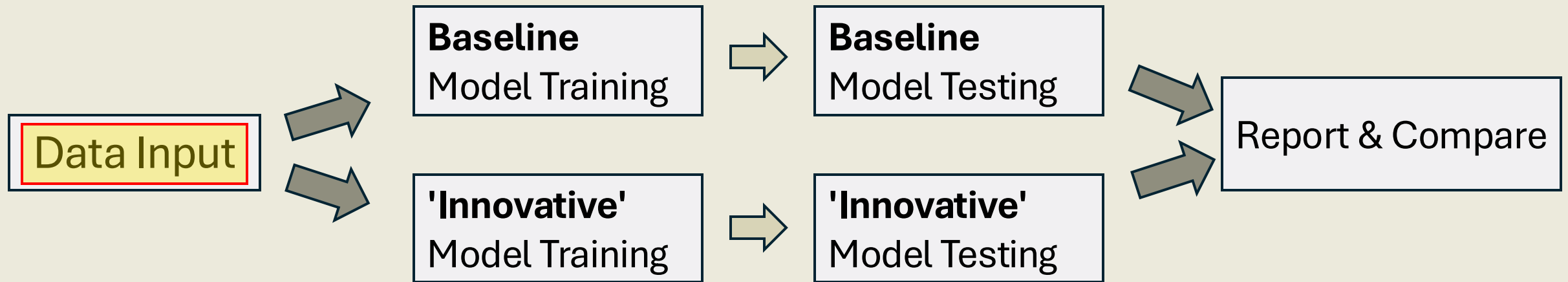
4.



How can we answer these questions?



Experimental Design



Our data input requires a wide range of needs in order to capture an entire scene with object detection capabilities.

Data Input

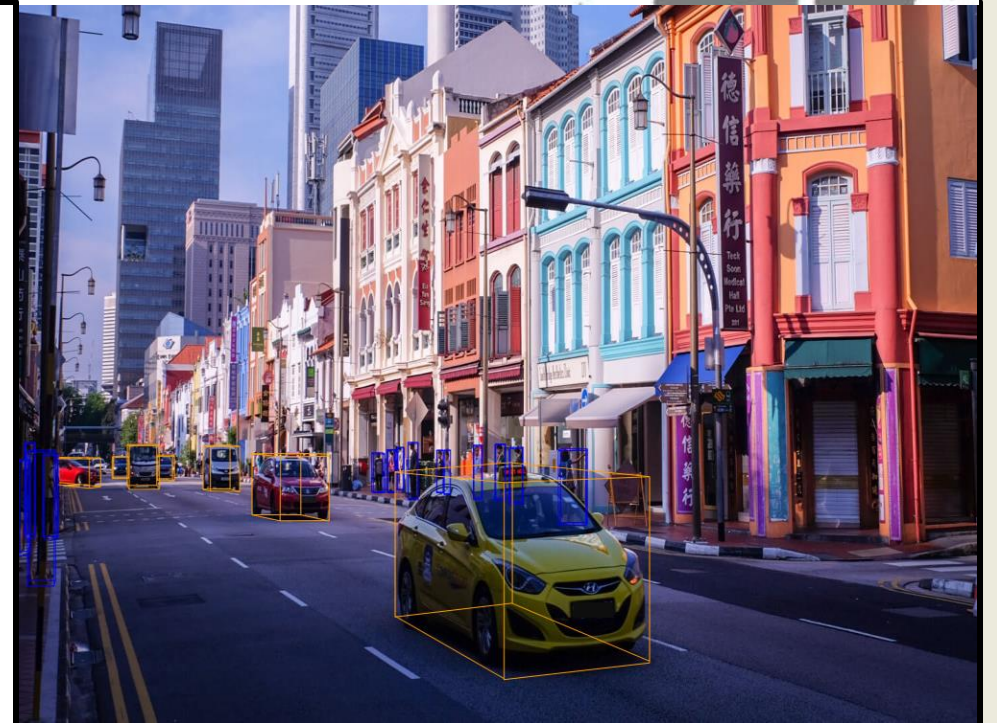
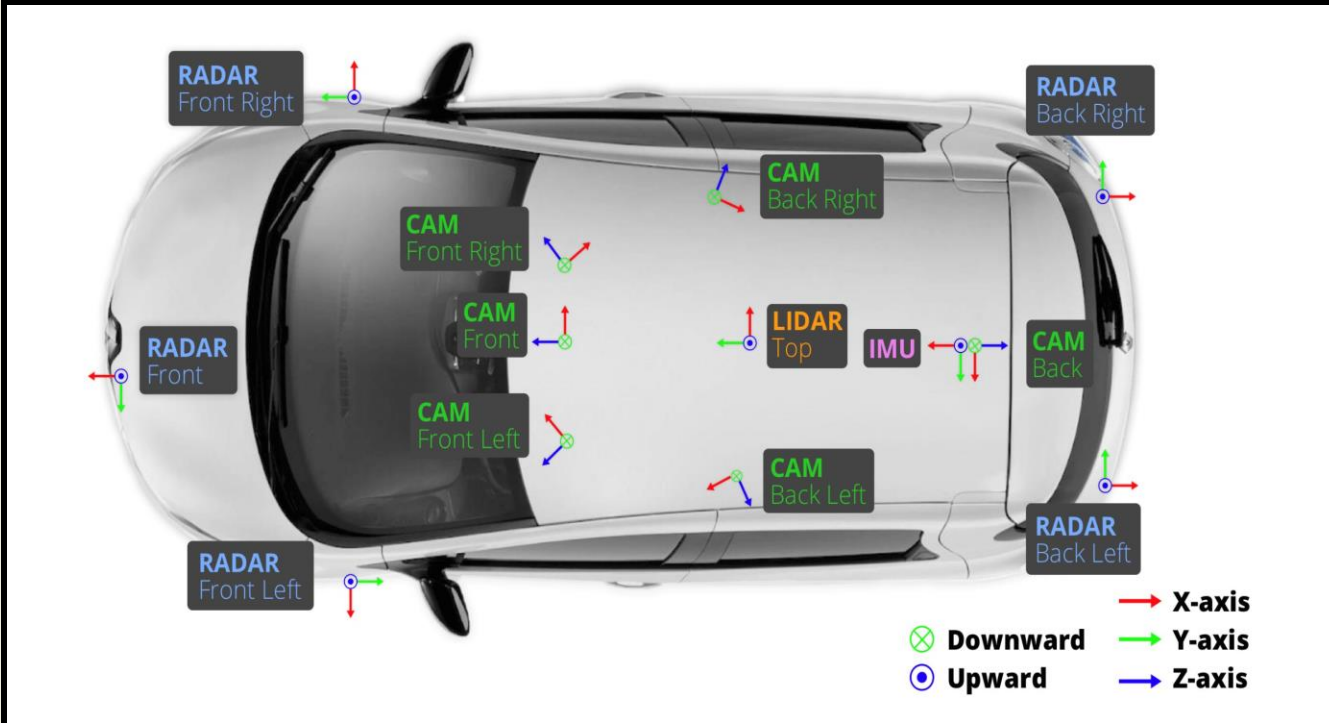
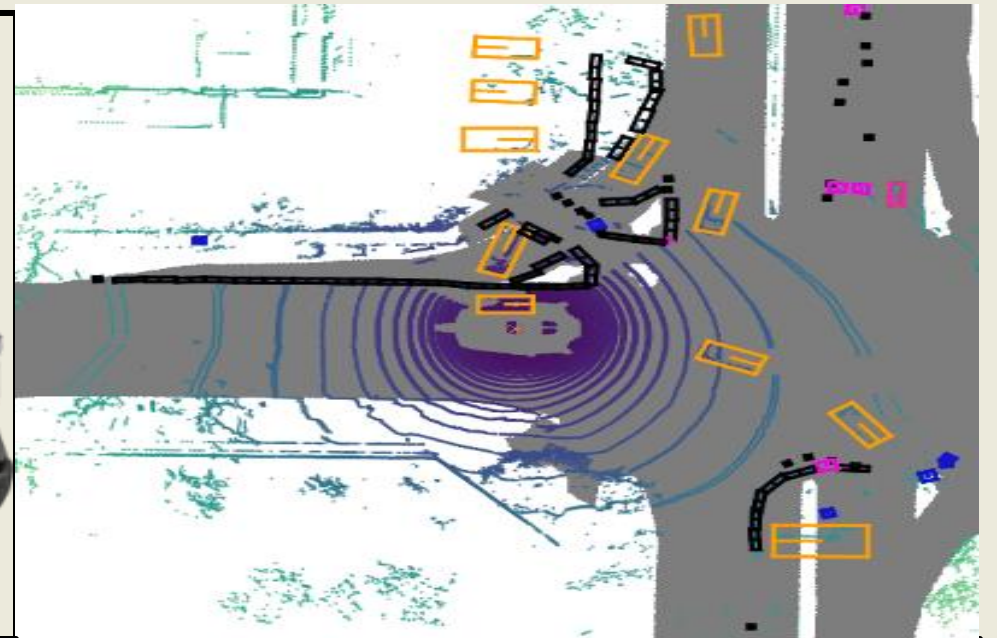
Continuous
images of a
scene

Labeled
Objects

Visibility of
Objects
Annotated

Intrinsic &
extrinsic
camera data

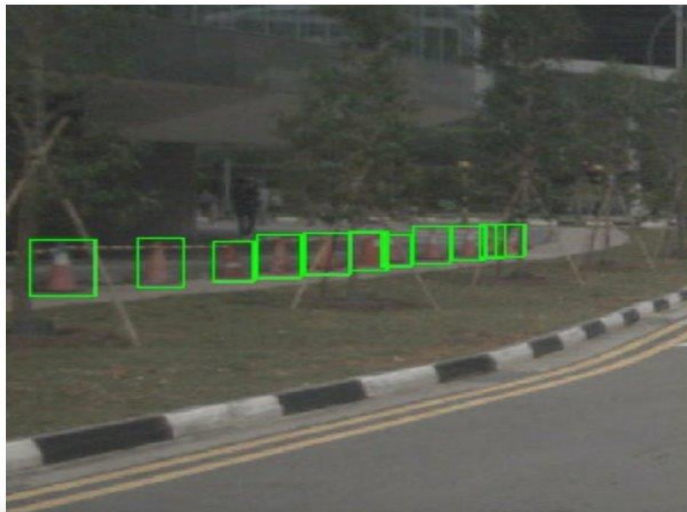
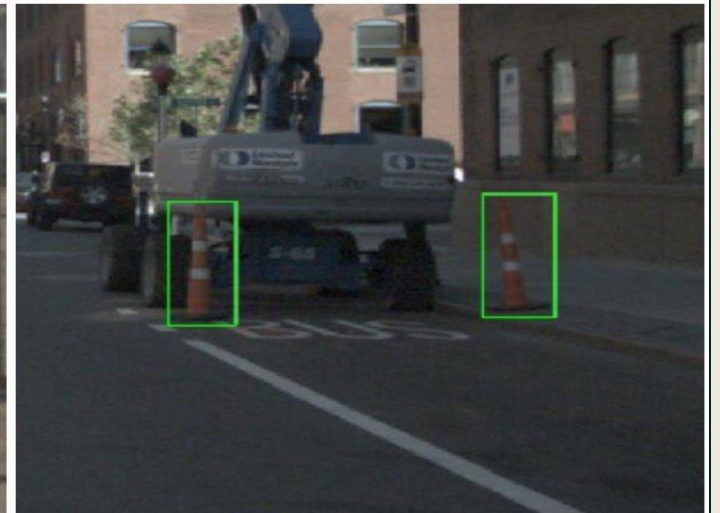
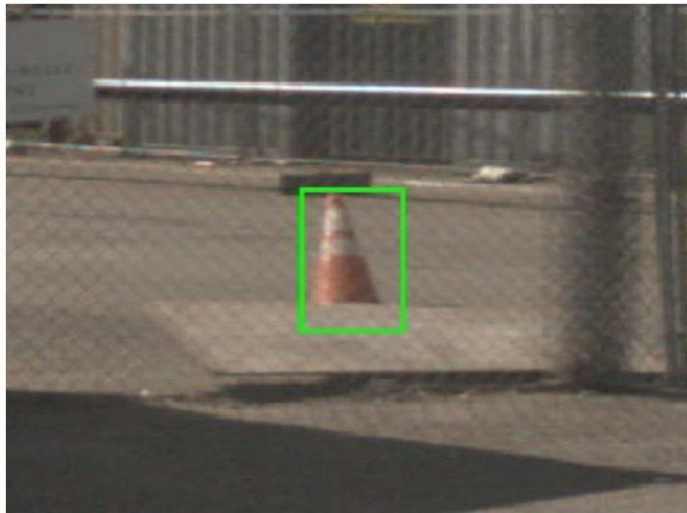
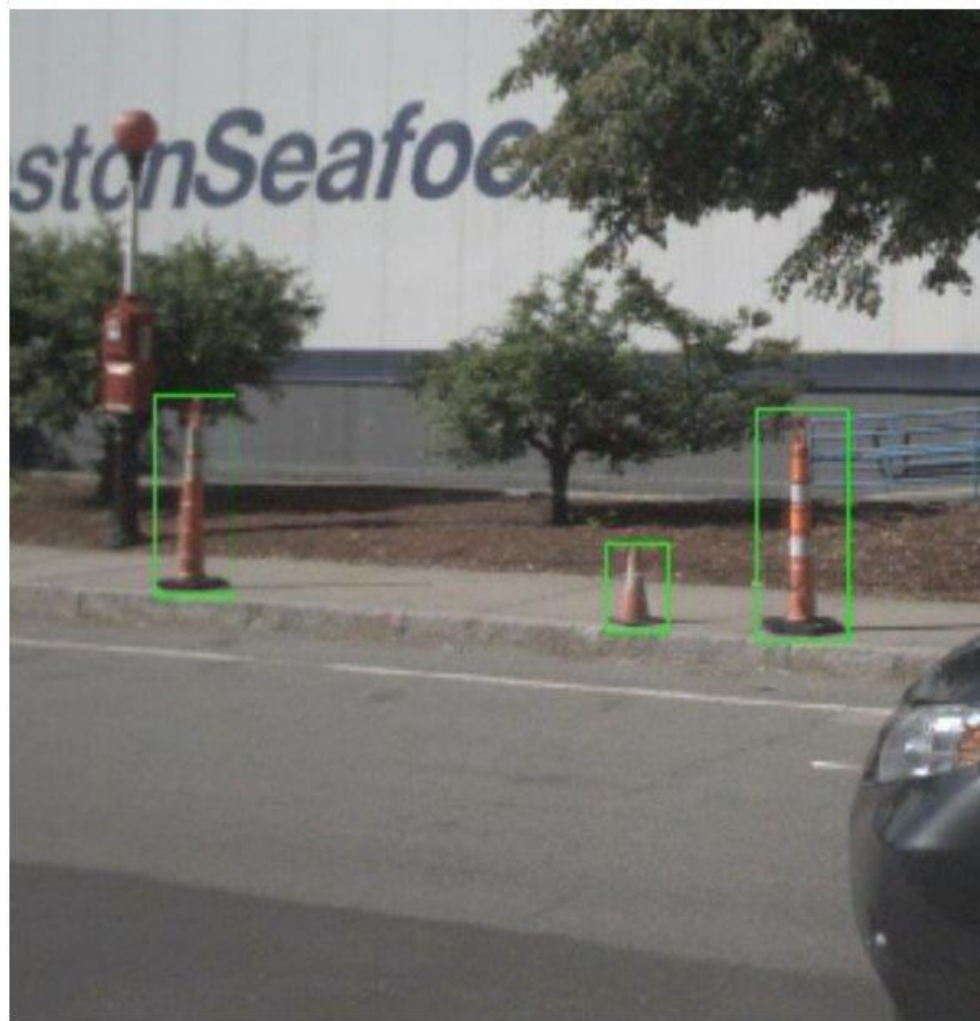
Adaptable
for ML
models



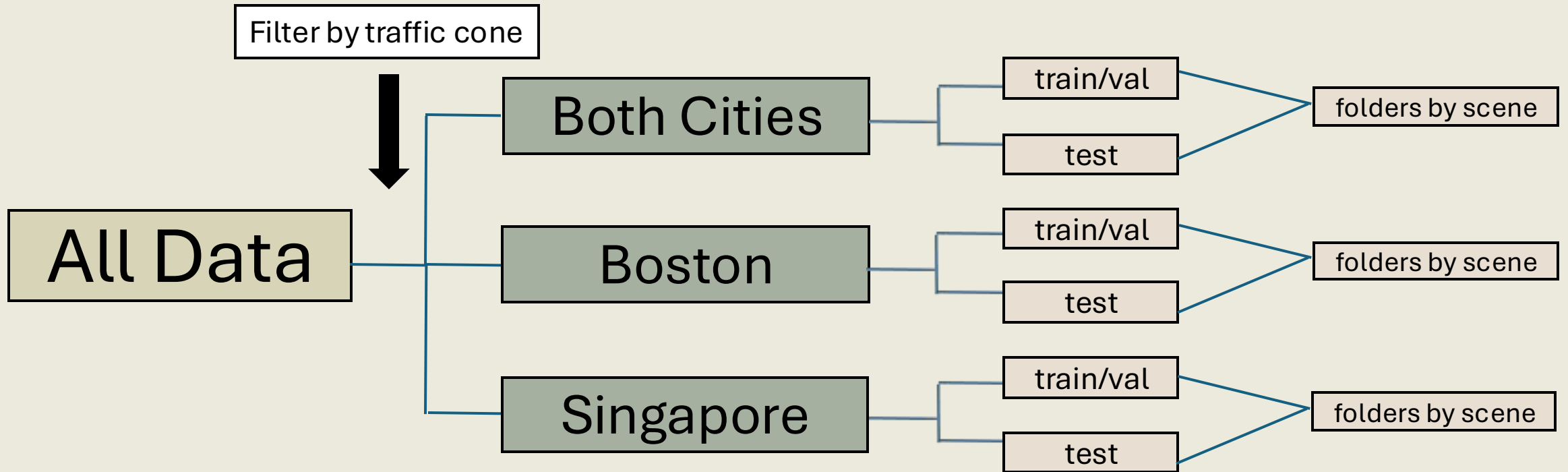


Traffic Cones

Source: https://github.com/nutonomy/nuscenes-devkit/blob/master/docs/instructions_nuscenes.md#traffic-cone



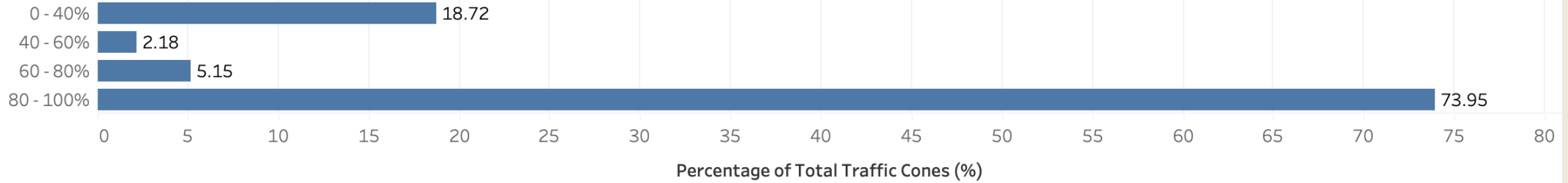
Data Structure



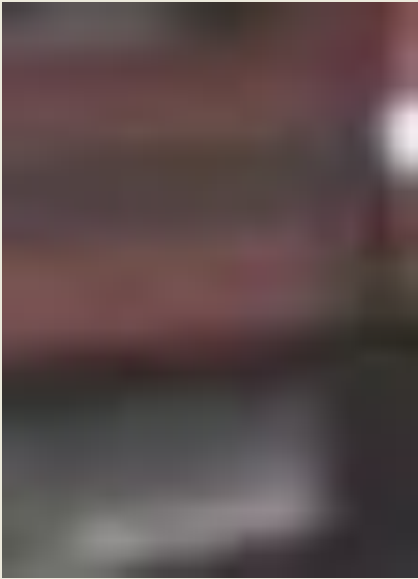
Nearly 74% of Traffic Cones Have High Visibility (80-100% Level)

Dataset sourced from NuScenes Mini Dataset--Total traffic cones analyzed: 1378

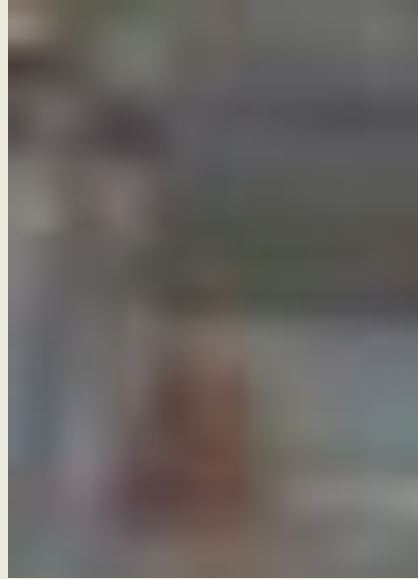
Visibility Level



Sum of Percentage for each Visibility Level. The marks are labeled by sum of Percentage.



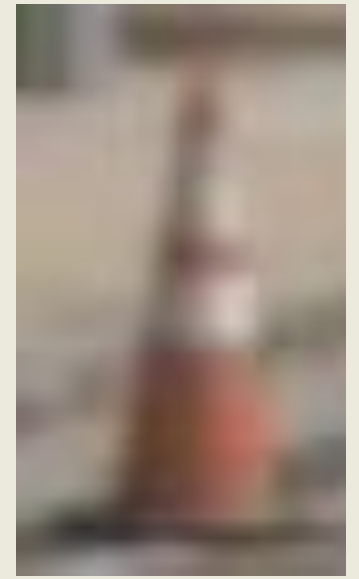
0-40%



40-60%

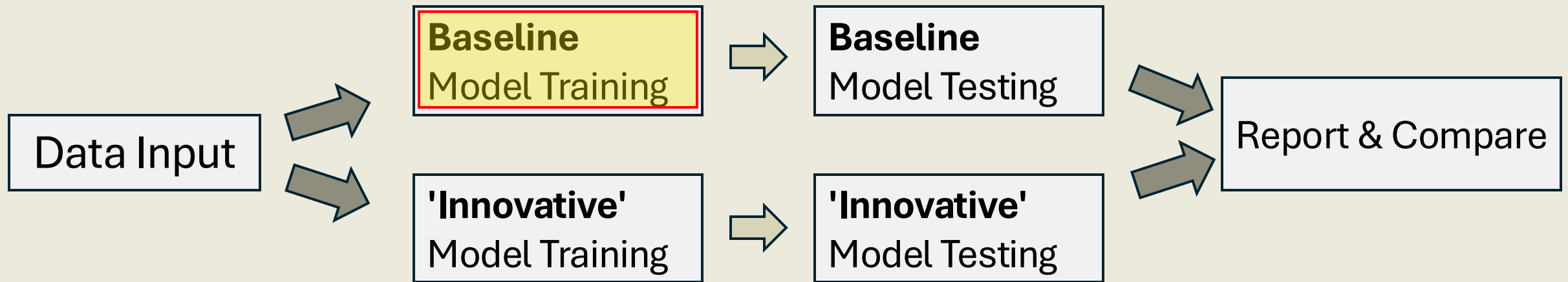


60-80%



80-100%

Experimental Design

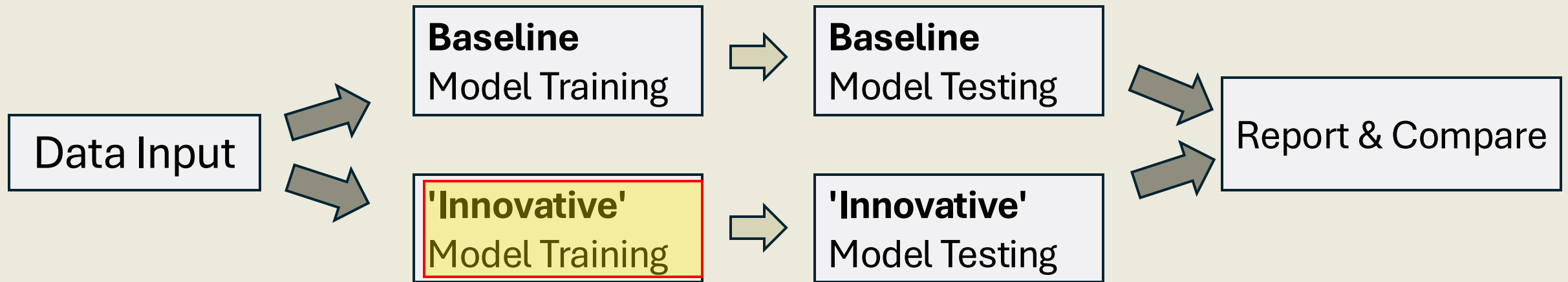


We define a baseline model as known reference, which would not consider obscurement or scene context between single images

You Only Look Once (YOLO)

- A real-time object detection algorithm
- Predict both the class and location of objects in an image
- Achieves high accuracy
- Fast object detection
- Good performance on small objects

Experimental Design



Our 'innovative' approach uses camera location to understand how the images connect to one another within a scene.

1. How much can partially obscured object detection improve using multi-views?	2. How do different levels of obscurement with multiple views affect detection accuracy?
3. What number of scene views do we need for high detection performance?	4. What features of our objects contribute to detected boxes?

DETR + 3D

Multi-view
image inputs

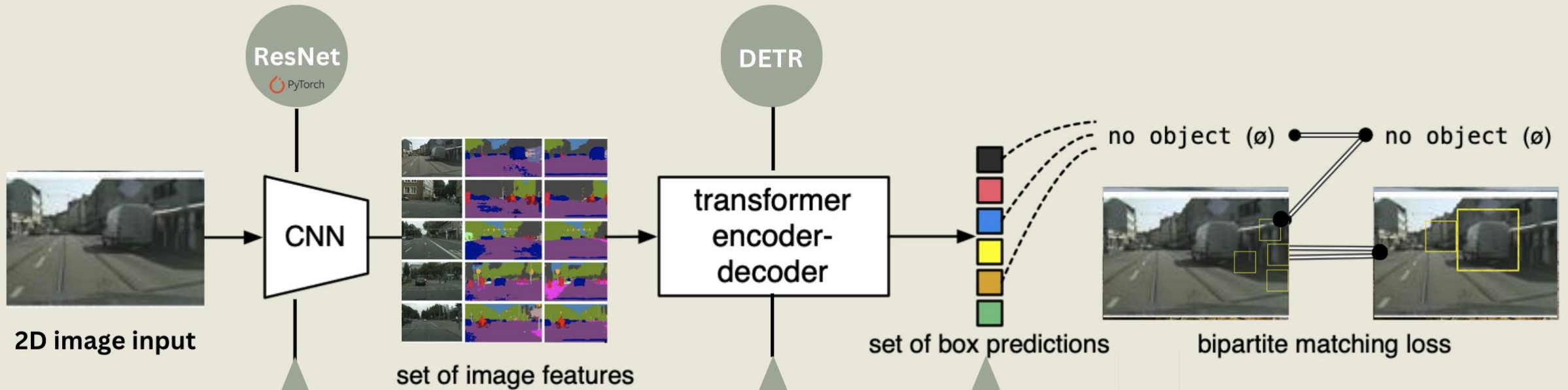
Obscured
object detection

Efficient
+ light-weight

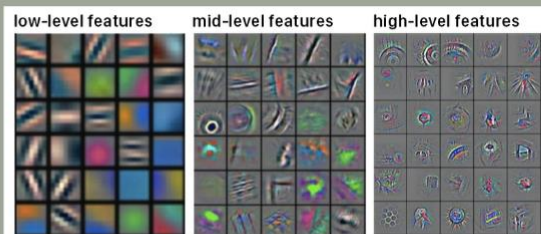
- Detects 3D objects directly from 2D multi-view images using transformer-based object queries

Method	NDS ↑	mAP ↑
Mono3D	0.429	0.366
DHNet	0.437	0.363
PGD [40]	0.448	0.386
DD3D [37] †	0.477	0.418
DETR3D (Ours) #	0.479	0.412

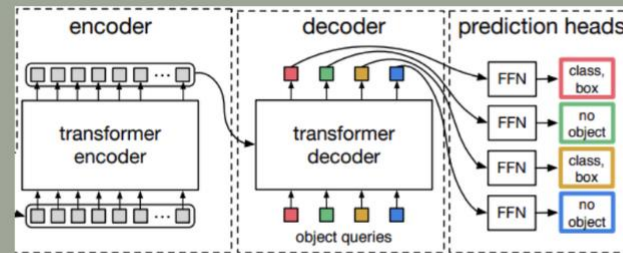
DETR + 3D



CNN Backbone:

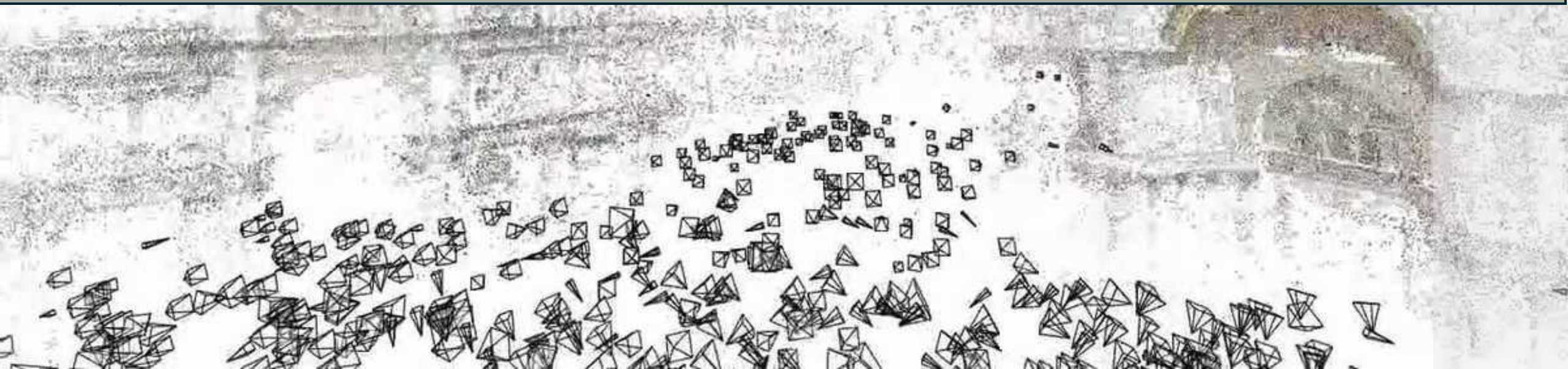


Transformer Architecture:





What's next?



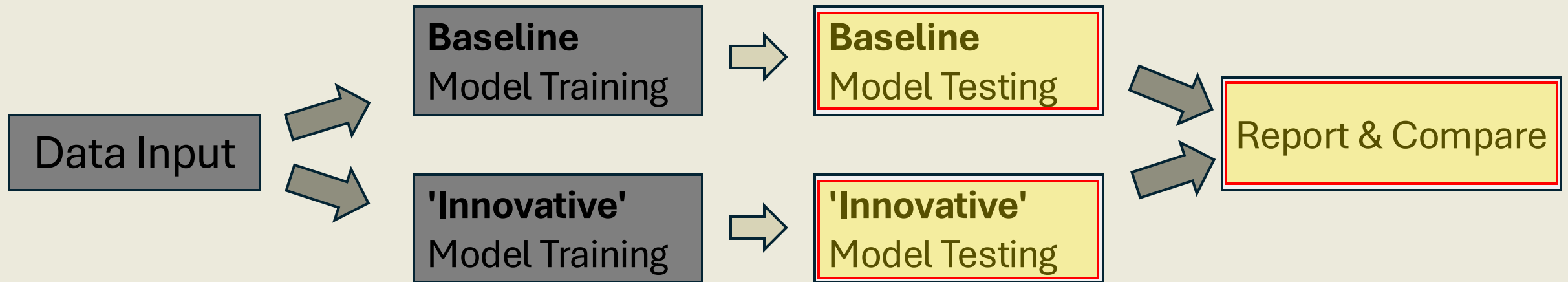
How much can partially
obscured object
detection improve using
1. multi-views?

How do different levels of
obscurement with
multiple views affect
detection accuracy? 2.

What number of scene
views do we need for
high detection
performance? 3.

What features of
our objects
contribute to
detected boxes? 4.

Spring Semester



We will test both our models and analyze the results so we can answer specifically the **interaction between multiple view and object visibility performance.**

TIMELINE

**Dec
2024**

Training procedures
complete

Testing of models and
evaluation metrics (Q2)

**Jan
2025**

**Feb
2025**

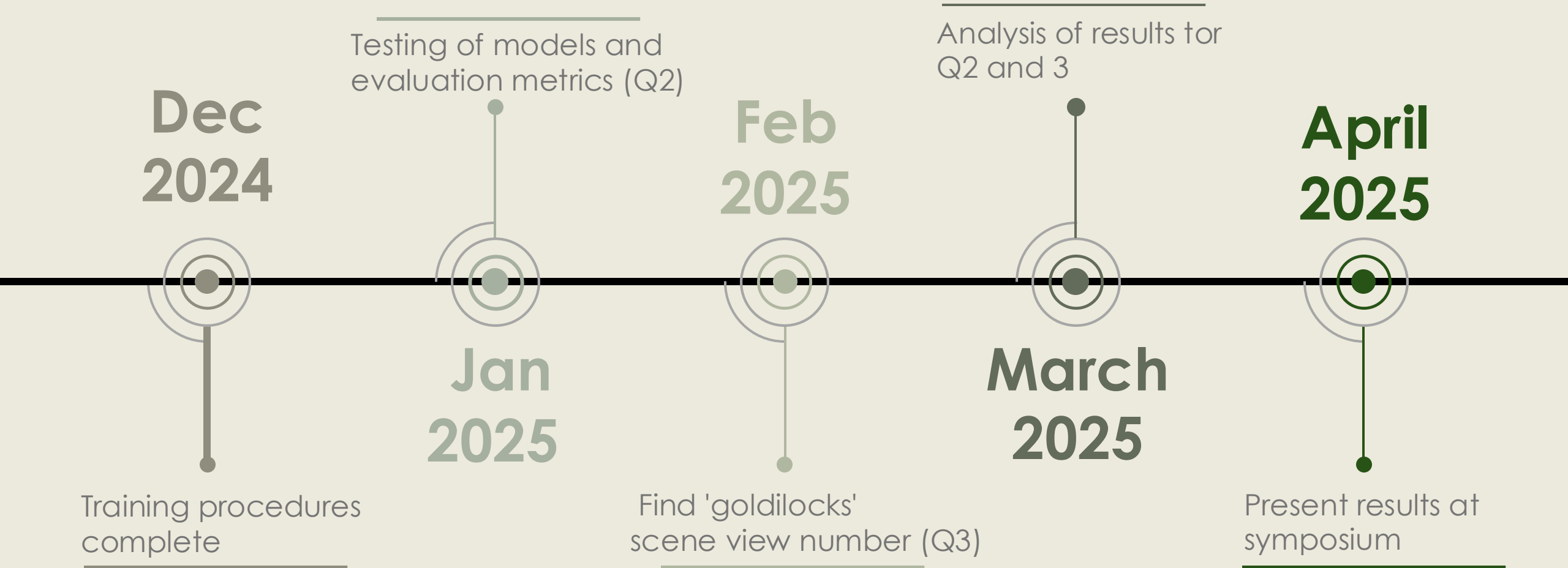
Find 'goldilocks'
scene view number (Q3)

Analysis of results for
Q2 and 3

**March
2025**

**April
2025**

Present results at
symposium



APPENDIX

DETR3D Processes and Description

