

Spatialyze: A Geospatial Video Analytics System with Spatial-Aware Optimizations

EPIC
DATA lab

Berkeley
UNIVERSITY OF CALIFORNIA

Overview

#1 Domain Challenges

- * Motivation
- * Usage Scenario
- * Typical Workflow

#3 System Optimization

- * Video Processor
- * Optimization Techniques

#2 System Interface

- * Data Model
- * Programming Model

#4 Evaluation

- * Experiment Setups
- * Results

Overview

#1 Domain Challenges

- * Motivation
- * Usage Scenario
- * Typical Workflow

#3 System Optimization

- * Video Processor
- * Optimization Techniques

#2 System Interface

- * Data Model
- * Programming Model

#4 Evaluation

- * Experiment Setups
- * Results

Overview

#1 Domain Challenges

- * Motivation
- * Usage Scenario
- * Typical Workflow

#3 System Optimization

- * Video Processor
- * Optimization Techniques

#2 System Interface

- * Data Model
- * Programming Model

#4 Evaluation

- * Experiment Setups
- * Results

Overview

#1 Domain Challenges

- * Motivation
- * Usage Scenario
- * Typical Workflow

#3 System Optimization

- * Video Processor
- * Optimization Techniques

#2 System Interface

- * Data Model
- * Programming Model

#4 Evaluation

- * Experiment Setups
- * Results

SPATIALYZE

#1 Domain Challenges

- * Motivation
- * Usage Scenario
- * Typical Workflow

#1 DOMAIN CHALLENGES

Motivation

What are Geospatial Videos?

#1 DOMAIN CHALLENGES

Motivation

Geospatial Videos



Autonomous Vehicle



Traffic Camera



Drone Camera

Motivation

Geospatial Videos



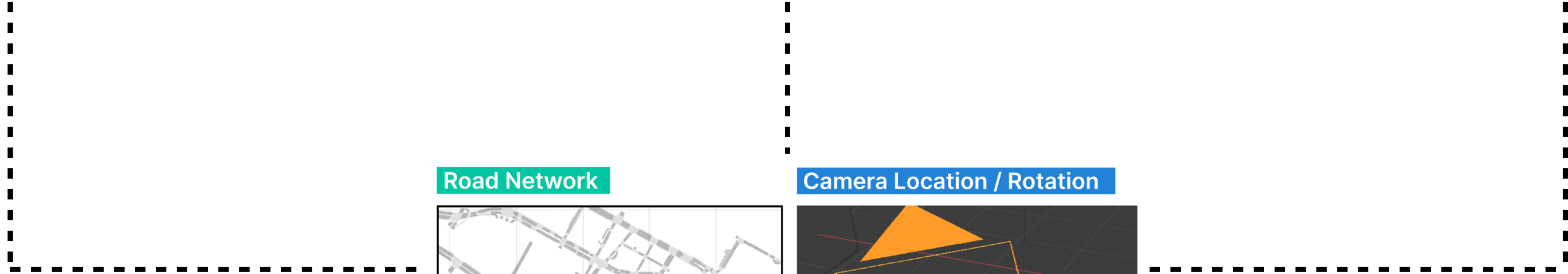
Autonomous Vehicle



Traffic Camera



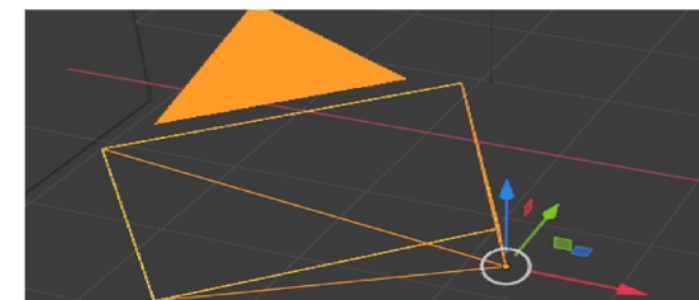
Drone Camera



Road Network



Camera Location / Rotation



Comes with geospatial metadata

Motivation



Autonomous Vehicle



Traffic Camera



Drone Camera

A data journalist investigates the behaviors of autonomous driving vehicles when each of them sees cars crash at an intersection.

Motivation



Autonomous Vehicle



Traffic Camera



Drone Camera

A traffic analyst wants to determine when is the busiest hours of a traffic intersection based on recorded traffic camera footages.

Motivation



Autonomous Vehicle



Traffic Camera



Drone Camera

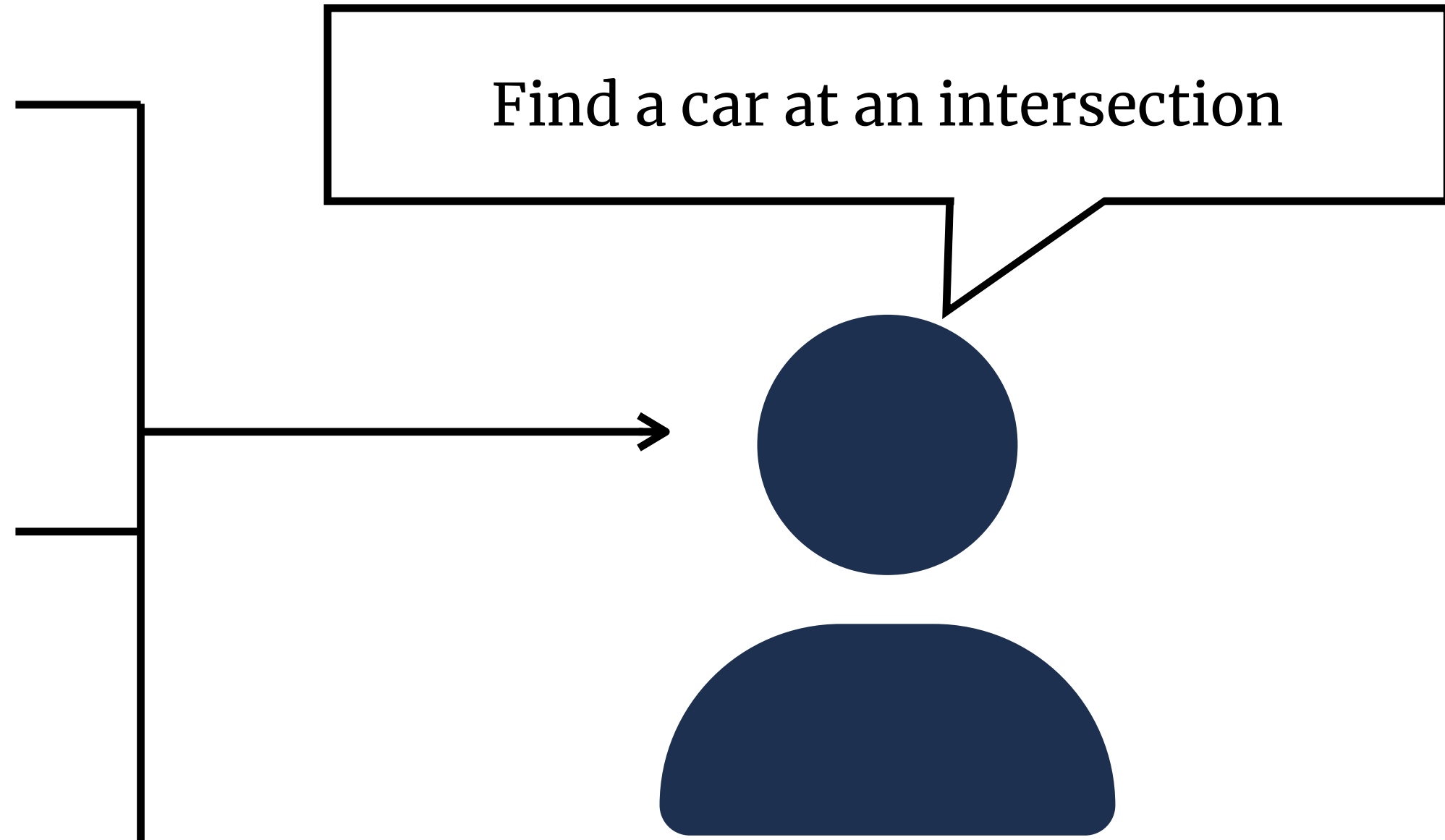
An engineer wants to develop an application that reports realtime parking spots availability from drone camera footages

#1 DOMAIN CHALLENGES

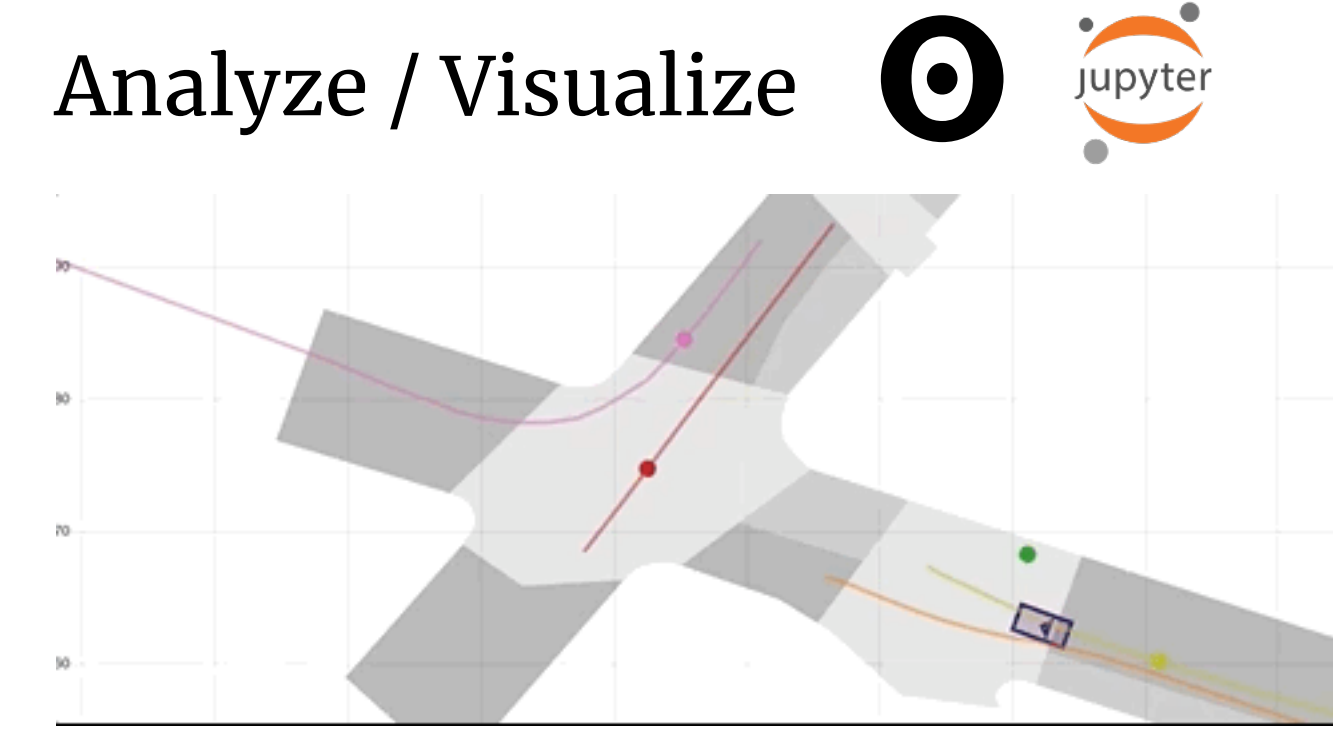
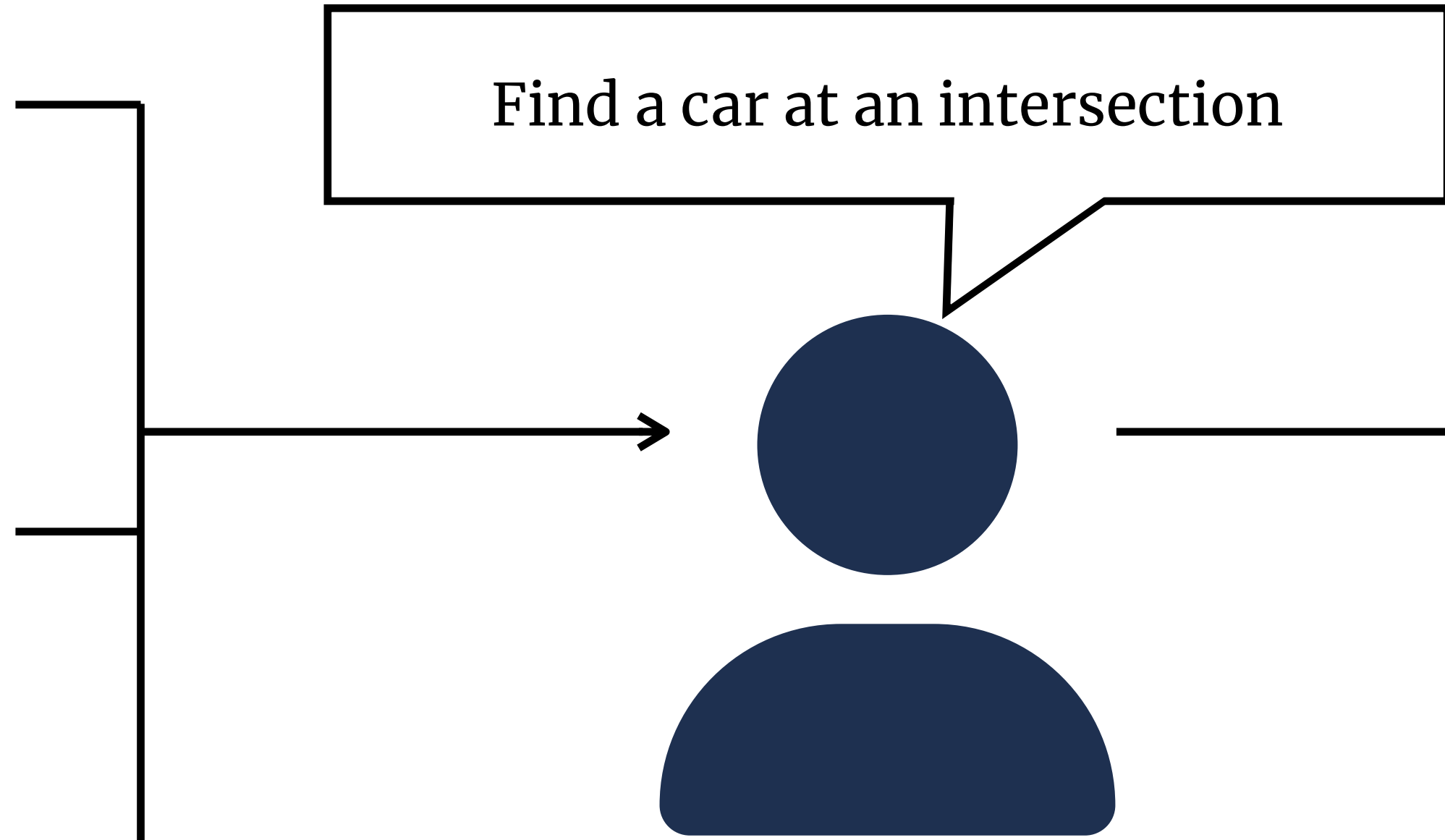
Usage Scenario



Usage Scenario



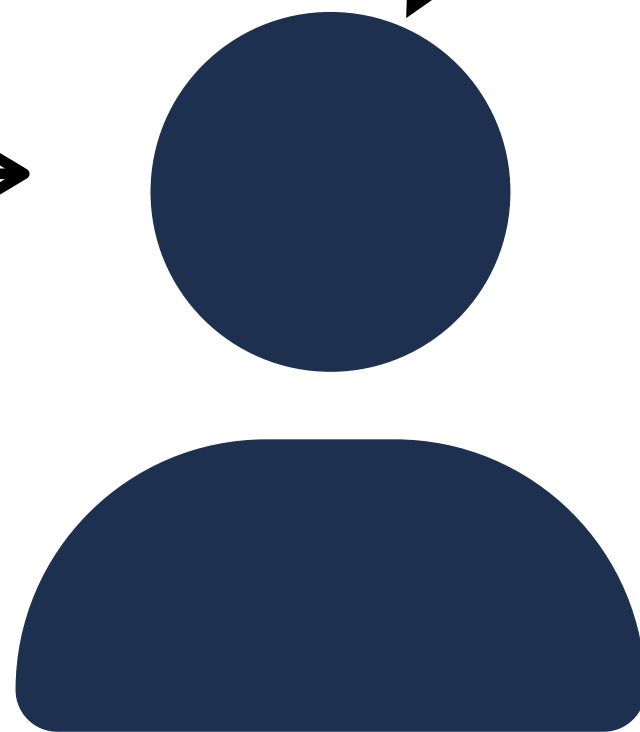
Usage Scenario



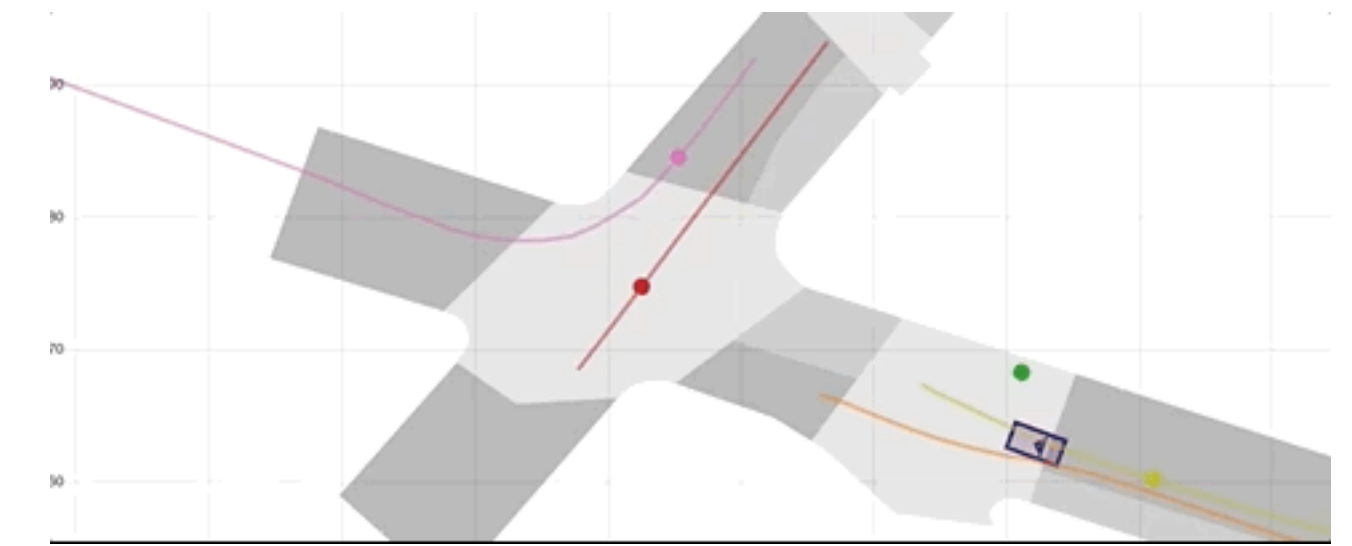
Usage Scenario



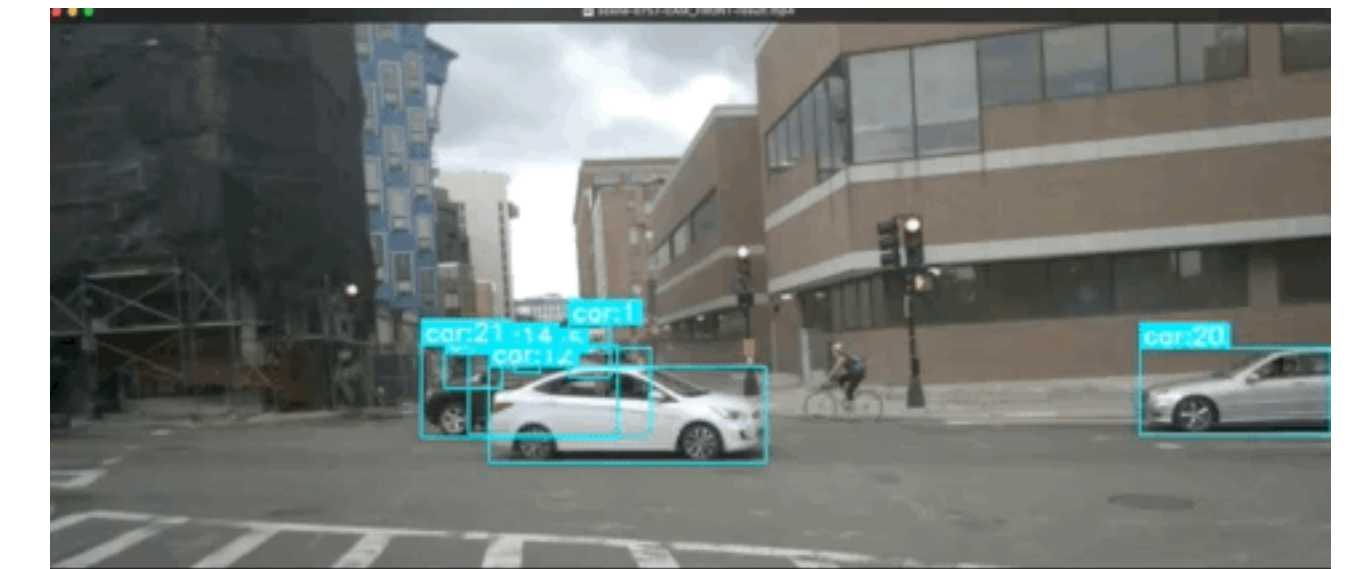
Find a car at an intersection



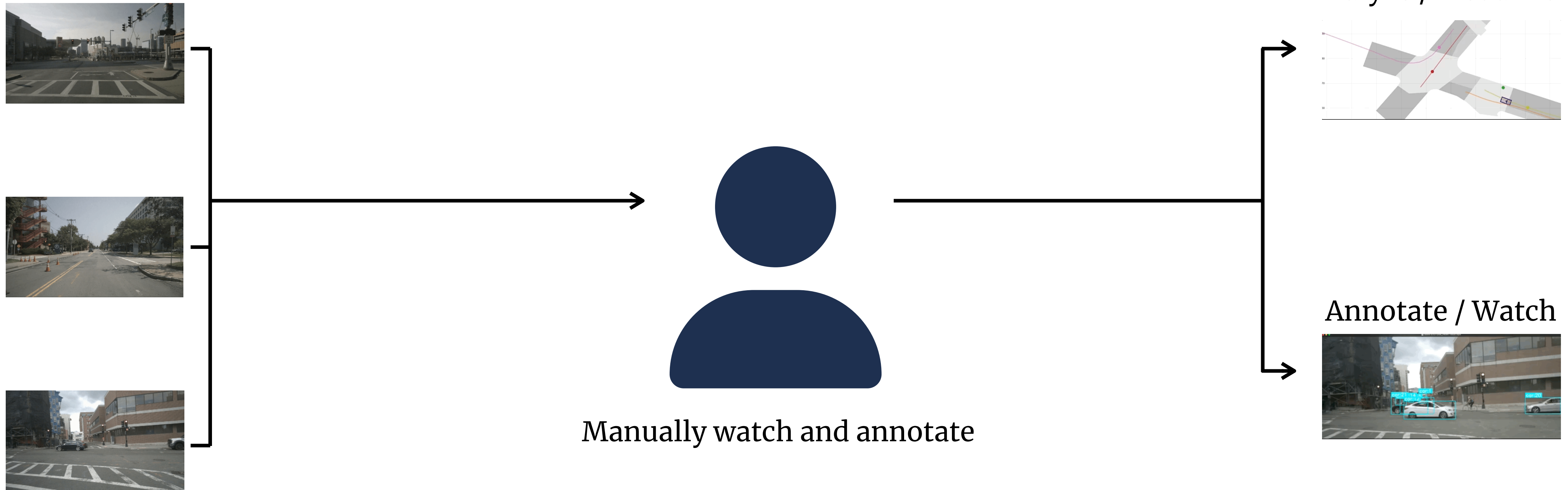
Analyze / Visualize



Annotate / Watch



Typical Workflow (Manual)



Typical Workflow (Manual)



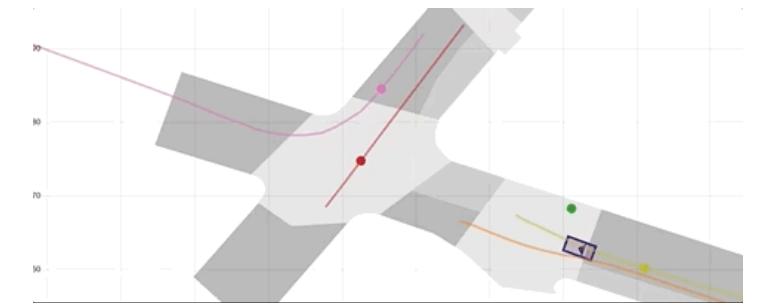
Manually watch and annotate



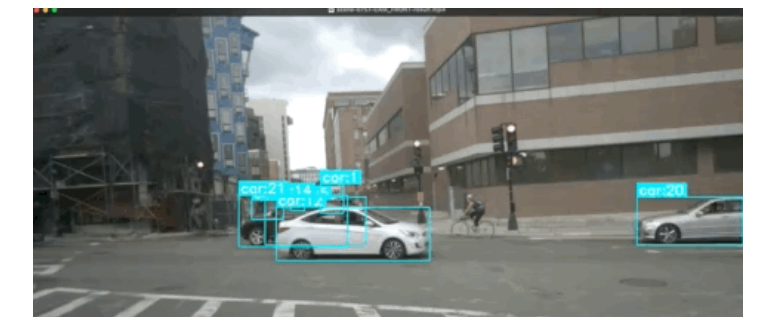
Time-consuming



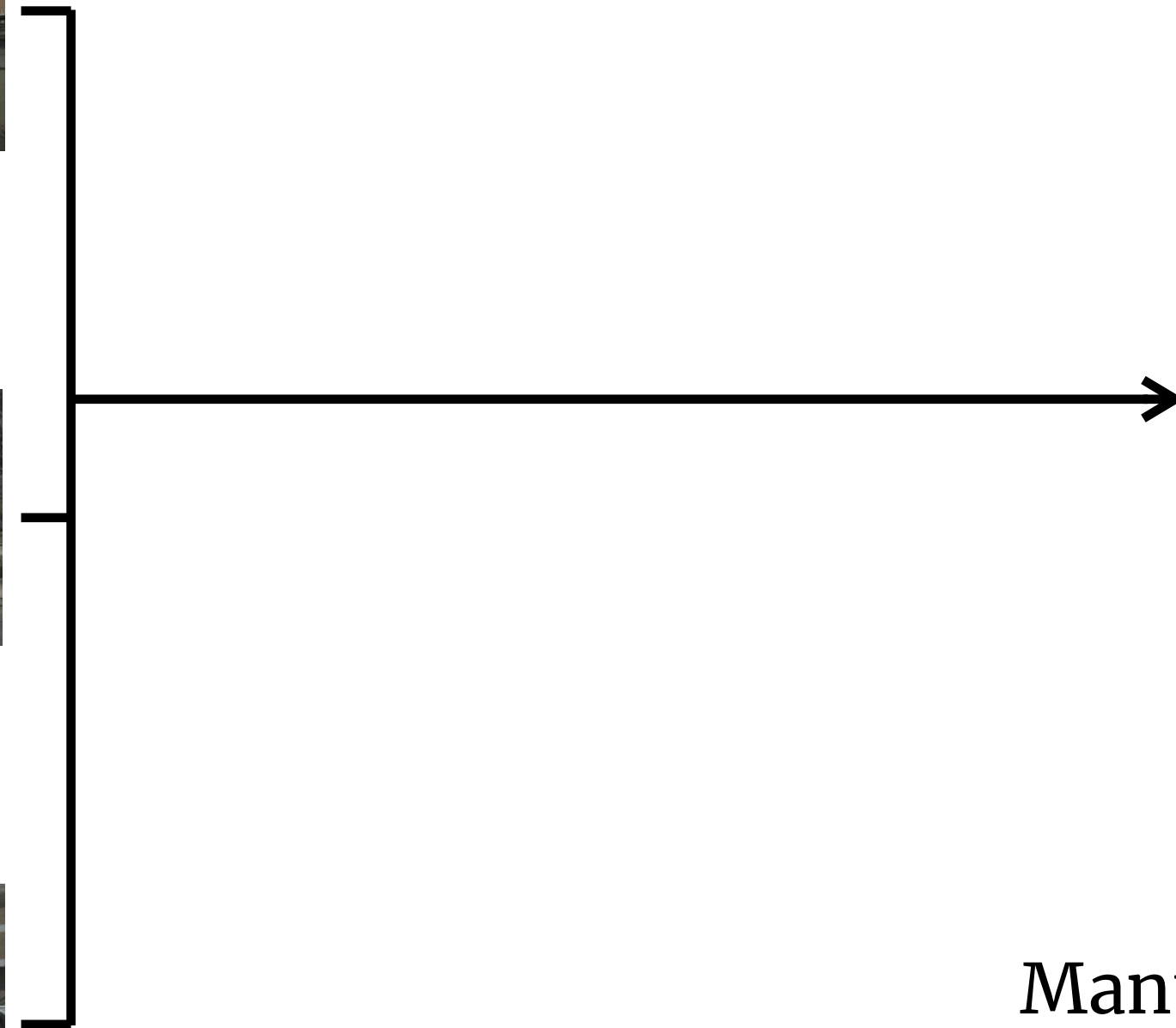
Analyze / Visualize



Annotate / Watch



Typical Workflow (Manual)



Manually watch and annotate



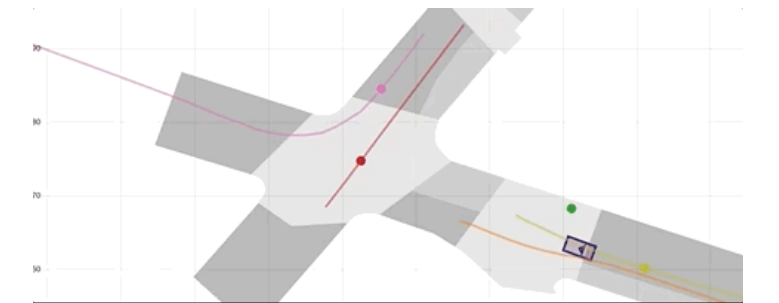
Time-consuming



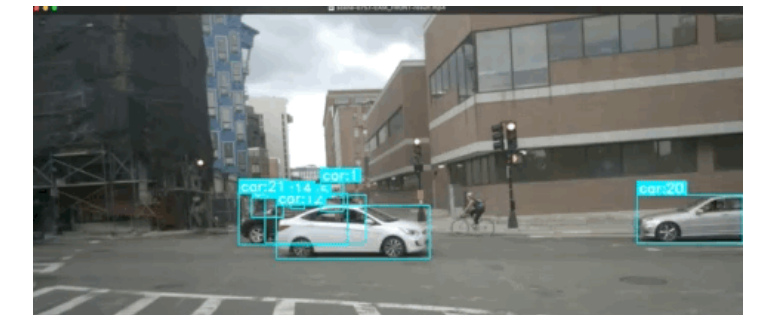
Tedious &
Error-prone



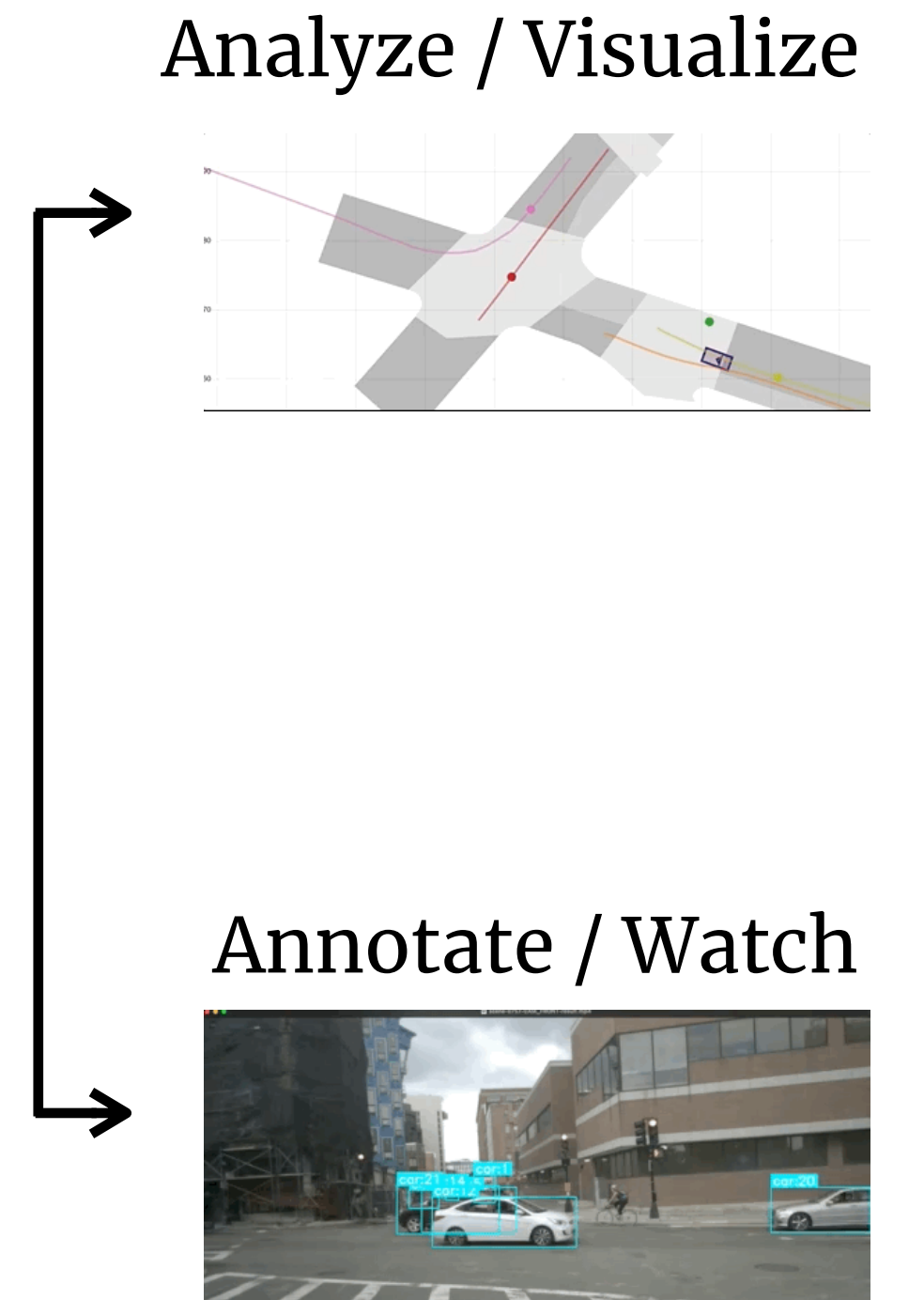
Analyze / Visualize



Annotate / Watch



Typical Workflow (Programming)



Typical Workflow (Programming)



Object Detector

```
with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: List[str] = class_mapping_to_list(_names)
    self.model.eval()

    assert isinstance(self.imgsz, List), type(self)
    self.model.warmup(imgsz=(1, 3, *self.imgsz)) #

    for frame_idx, im0 in enumerate(self.frames.stream):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _ = letterbox(im0, self.imgsz, stride=self.stride, auto=True)
        im = im.transpose(2, 0, 1)[::-1] # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

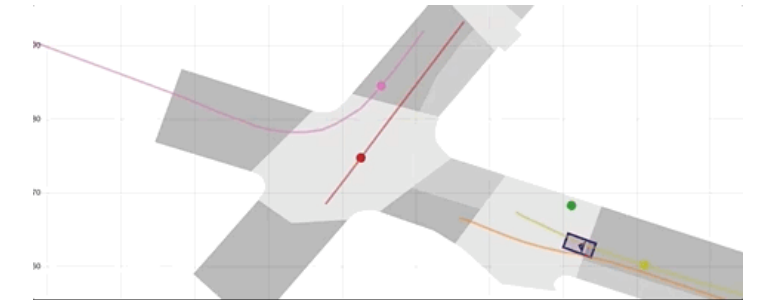
        # t1 = time_sync()
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        im /= 255.0 # 0 - 255 to 0.0 - 1.0
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        # Inference
        pred = self.model(im, augment=self.augment)

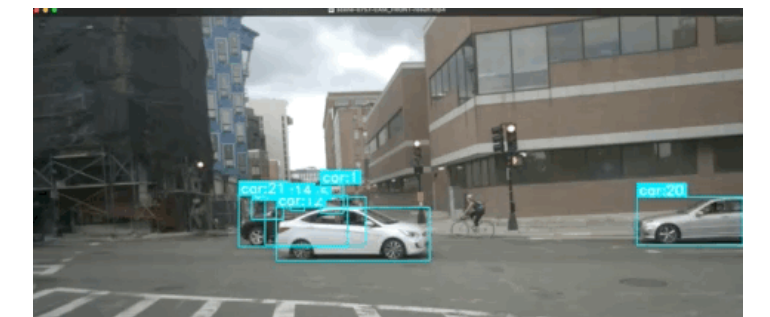
        # Apply NMS
        pred = non_max_suppression(
            pred,
            self.conf_thres,
            self.iou_thres,
            self.classes,
            self.agnostic_nms,
            max_det=self.max_det,
        )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4])
        yield Detection2D(
            det, names, [DetectionId(frame_idx, ord)]
        )
```

Analyze / Visualize



Annotate / Watch



Typical Workflow (Programming)



Object Detector

```
with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(self.model.eval())

    assert isinstance(self.imsz, List), type(self.imsz)
    self.model.warmup(imsz=[1, 3, *self.imsz])

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imsz, stride=self.stride)
        im = im.transpose(2, 0, 1)[::-1] # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        im /= 255.0 # 0 - 255 to 0.0 - 1.0
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        # Inference
        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det)

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4])
        yield Detection2D(det, names, [DetectionId(frame_idx, ord)])
```

3D Location Estimator

```
class LocationEstimator:
    def __init__(self, model, device, names, metric_callback=None):
        self.model = model
        self.device = device
        self.names = names
        self.metric_callback = metric_callback

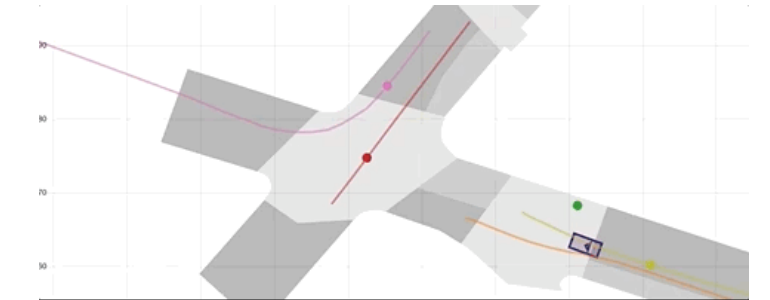
    def estimate(self, frame_idx, im0):
        # Preprocess
        im, _, _ = letterbox(im0, self.imsz, stride=self.stride)
        im = im.transpose(2, 0, 1)[::-1] # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        im /= 255.0 # 0 - 255 to 0.0 - 1.0
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

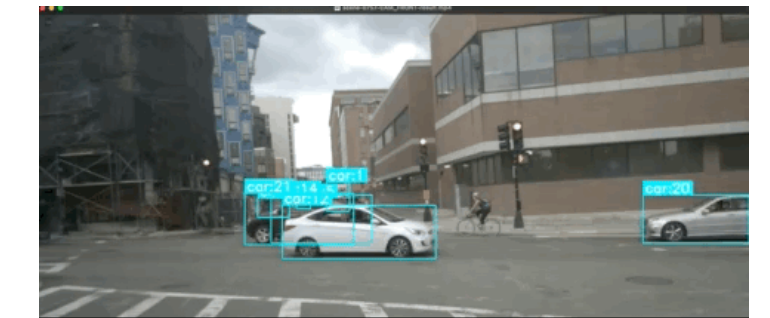
        # Inference
        pred = self.model(im, augment=self.augment)

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4])
        yield Detection2D(det, names, [DetectionId(frame_idx, ord)])
```

Analyze / Visualize



Annotate / Watch



Typical Workflow (Programming)



```
Object Detector

with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(
        self.model.eval())

    assert isinstance(self.imsz, List), type(self.
        self.model.warmup(imsz=(1, 3, *self.imsz)) #

    for frame_idx, im0 in enumerate(self.frames.str
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imsz, stride
        im = im.transpose(2, 0, 1)[::-1] # HWC t
        im = np.ascontiguousarray(im) # contiguous

        # t1 = time_sync()
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        # Inference
        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(
            pred,
            self.conf_thres,
            self.iou_thres,
            self.classes,
            self.agnostic_nms,
            max_det=self.max_det,
        )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(
            det), :4] = scale_boxes(im.shape[2:], det)
        yield Detection2D(
            det, names, [DetectionID(frame_idx, ord
```

```
3D Location Estimator

class BaseEstimator:
    """Base class for 3D location estimators.

    Attributes:
        img_size: tuple[int, int]
        stride: int
        device: torch.device
        names: list[str]
        conf_thres: float
        iou_thres: float
        classes: list[int]
        agnostic_nms: bool
        max_det: int
    """

    def __init__(self, img_size: tuple[int, int], stride: int, device: torch.device, names: list[str], conf_thres: float, iou_thres: float, classes: list[int], agnostic_nms: bool, max_det: int):
        self.img_size = img_size
        self.stride = stride
        self.device = device
        self.names = names
        self.conf_thres = conf_thres
        self.iou_thres = iou_thres
        self.classes = classes
        self.agnostic_nms = agnostic_nms
        self.max_det = max_det

    def process(self, frame_idx: int, im0: torch.Tensor) -> list[Detection3D]:
        """Process a frame to estimate 3D locations.

        Args:
            frame_idx: int
            im0: torch.Tensor

        Returns:
            list[Detection3D]
        """
        # ... (implementation details) ...
        return detections
```

```
Object Tracker

save_detections: list[list[dict, torch.Tensor]] = []
class ListList | None = None
def __init__(self) -> None:
    """Initialize the tracker.

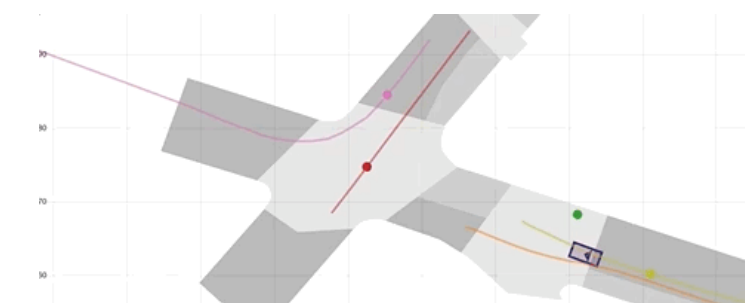
    Attributes:
        save_detections: list[list[dict, torch.Tensor]]
        class_mapping: dict[str, int]
        camera_configs: list[CameraConfig]
        tracks: list[Track]
        detections: list[list[dict, torch.Tensor]]
        camera_configs: list[CameraConfig]
        track_id: int
        track_id_map: dict[int, int]
        track_id_counter: int
    """
    self.save_detections = save_detections
    self.class_mapping = class_mapping
    self.camera_configs = camera_configs
    self.tracks = tracks
    self.detections = detections
    self.camera_configs = camera_configs
    self.track_id = track_id
    self.track_id_map = track_id_map
    self.track_id_counter = track_id_counter

    def track(self, frame_idx: int, detections: list[Detection2D], camera_config: CameraConfig) -> list[Track]:
        """Track objects in a frame.

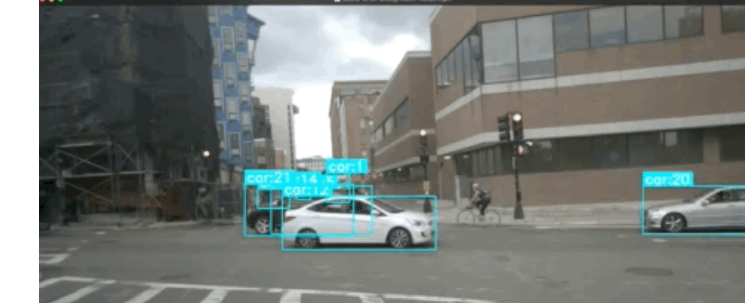
        Args:
            frame_idx: int
            detections: list[Detection2D]
            camera_config: CameraConfig

        Returns:
            list[Track]
        """
        # ... (implementation details) ...
        return tracks
```

Analyze / Visualize



Annotate / Watch



Typical Workflow (Programming)



Object Detector

```
with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(
        self.model.eval())

    assert isinstance(self.imgsz, List), type(self)
    self.model.warmup(imgsz=(1, 3, *self.imgsz))

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imgsz, stride=
            im = im.transpose((2, 0, 1)) # HWC to CHW
            im = np.ascontiguousarray(im) # contiguous

        # Inference
        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(
            pred,
            self.conf_thres,
            self.iou_thres,
            self.classes,
            self.agnostic_nms,
            max_det=self.max_det,
        )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det,
            det, names, [DetectionID(frame_idx, ord)
```

3D Location Estimator

```
class BaseEstimator:
    def __init__(self, imgsz: tuple[int, int, int], device: str):
        self.imgsz = imgsz
        self.device = device

    def estimate(self, im: torch.Tensor) -> List[Detection3D]:
        # Preprocess image
        im = letterbox(im, self.imgsz, stride=
        im = im.transpose((2, 0, 1)) # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(
            pred,
            self.conf_thres,
            self.iou_thres,
            self.classes,
            self.agnostic_nms,
            max_det=self.max_det,
        )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det,
            det, names, [DetectionID(frame_idx, ord)
```

Object Tracker

```
save_detections: list[list, torch.Tensor] = []
class Tracker:
    def __init__(self, imgsz: tuple[int, int, int], device: str):
        self.imgsz = imgsz
        self.device = device

    def track(self, im: torch.Tensor) -> List[Track]:
        # Preprocess image
        im = letterbox(im, self.imgsz, stride=
        im = im.transpose((2, 0, 1)) # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        pred = self.model(im, augment=self.augment)

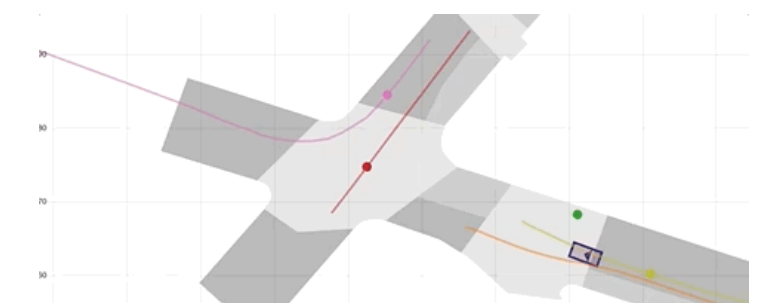
        # Apply NMS
        pred = non_max_suppression(
            pred,
            self.conf_thres,
            self.iou_thres,
            self.classes,
            self.agnostic_nms,
            max_det=self.max_det,
        )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det,
            det, names, [DetectionID(frame_idx, ord)
```

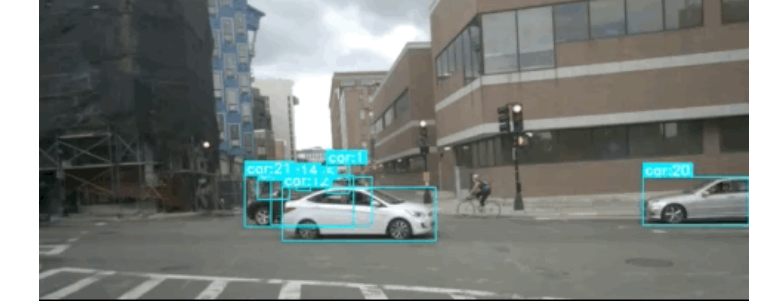
Filter Objects

```
select c0.frameNum,
        c0.cameraId,
        c0.filename,
        t0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
    SegmentPolygon.polygon,
    valueAtTimestamp(
        t0.translations,
        c0.timestamp
    )
)
```

Analyze / Visualize



Annotate / Watch

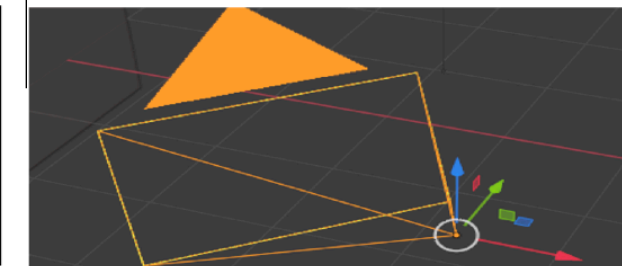


Typical Workflow (Pro)

Road Network



Camera Location / Rotation



Object Detector

```
with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(self.model.eval())

    assert isinstance(self.imgsz, List), type(self)
    self.model.warmup(imgsz=(1, 3, *self.imgsz)) #

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imgsz, stride=self.stride, pad=0) # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        # Apply NMS
        pred = self.model(im, augment=self.augment)
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det, )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4], im.shape)
        yield Detection2D(det, names, [DetectionID(frame_idx, ord)])
```

3D Location Estimator

```
class LocationEstimator:
    def __init__(self, road_network, camera_location, camera_rotation):
        self.road_network = road_network
        self.camera_location = camera_location
        self.camera_rotation = camera_rotation

    def estimate_location(self, frame_idx, detection):
        # ... (omitted code) ...
```

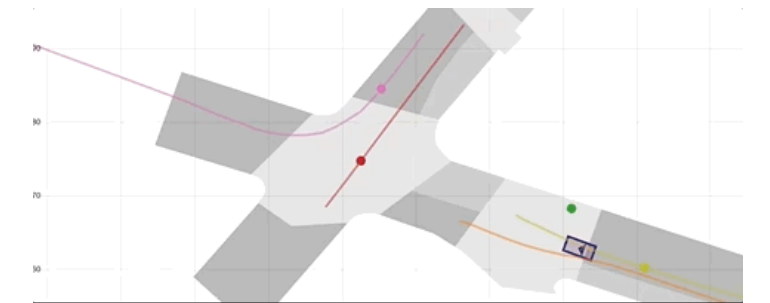
Object Tracker

```
def track(self, frame_idx, detection):
    # ... (omitted code) ...
```

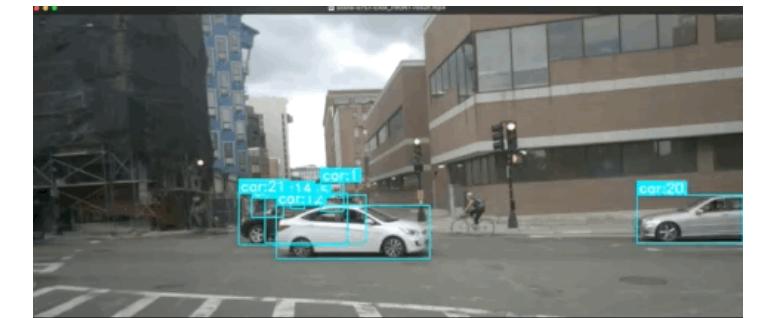
Filter Objects

```
select c0.frameNum,
       c0.cameraId,
       c0.filename,
       t0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
    SegmentPolygon.polygon,
    valueAtTimestamp(
        t0.translations,
        c0.timestamp
    )
)
```

Analyze / Visualize

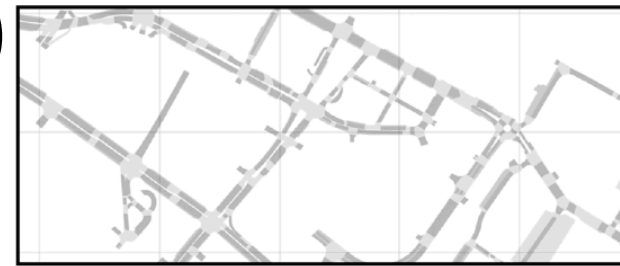


Annotate / Watch

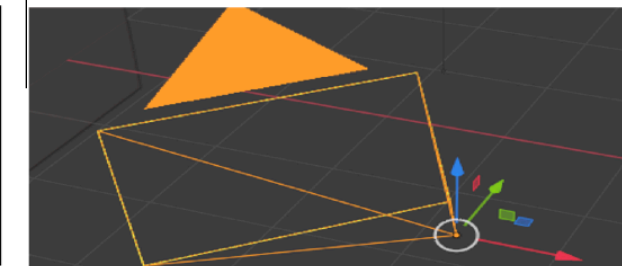


Typical Workflow (Pro)

Road Network



Camera Location / Rotation



Object Detector

```
with torch.no_grad():
    _names = self.model.names
    names: list[str] = class_mapping_to_list_names(self.model.eval())

    assert isinstance(self.imgsz, List), type(self.model.warmup(imgsz=(1, 3, *self.imgsz))) #

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imgsz, stride=self.stride, auto=True) # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        # Apply NMS
        pred = self.model(im, augment=self.augment)
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det, )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4], im.shape) # to original image coordinates
        yield Detection2D(
            det, names, [DetectionID(frame_idx, ord)]
```

3D Location Estimator

```
class LocationEstimator:
    def __init__(self, road_network, camera_config):
        self.road_network = road_network
        self.camera_config = camera_config

    def estimate_location(self, detections):
        # ... (omitted code) ...
        return location
```

Object Tracker

```
class ObjectTracker:
    def __init__(self, road_network, camera_config):
        self.road_network = road_network
        self.camera_config = camera_config

    def track_objects(self, detections):
        # ... (omitted code) ...
        return tracks
```

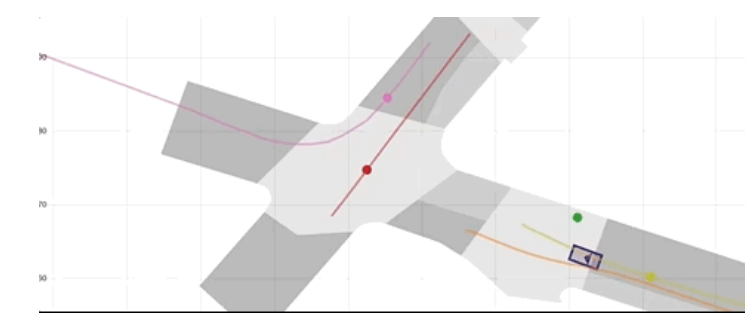
Filter Objects

```
select c0.frameNum,
       c0.cameraId,
       c0.filename,
       t0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
    SegmentPolygon.polygon,
    valueAtTimestamp(
        t0.translations,
        c0.timestamp
    )
)
```

Format Outputs

```
def format_outputs(tracks, detections, camera_configs):
    # ... (omitted code) ...
    return formatted_tracks
```

Analyze / Visualize

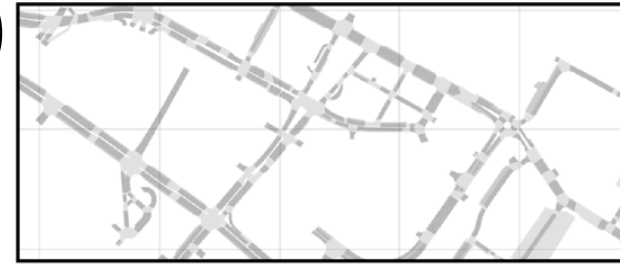


Annotate / Watch

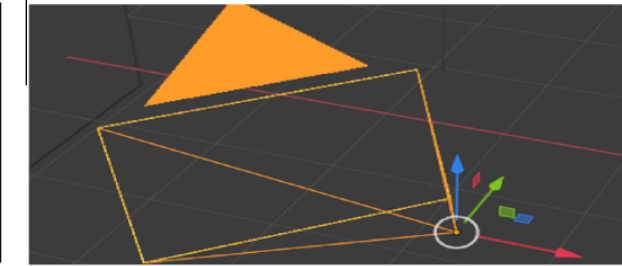


Typical Workflow (Pro)

Road Network



Camera Location / Rotation



Object Detector

```

with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(self.model.eval())

    assert isinstance(self.imgsz, List), type(self.model.warmup(imgsz=(1, 3, *self.imgsz))) #

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imgsz, stride=self.stride, auto=True) # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # t1 = time_sync()
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        # Inference
        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det, )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4], im0.shape) # xyxy
        yield Detection2D(det, names, [DetectionFrame, ...])
  
```

3D Location Estimator

```

class LocationEstimator:
    def __init__(self, road_network, camera_location):
        self.road_network = road_network
        self.camera_location = camera_location

    def estimate_location(self, det):
        # ... (code for estimating 3D location based on road network and camera view)
  
```

Object Tracker

```

save_detections: list[list, torch.Tensor] = []
class ListList: None = None
str: str = None

for detection, im0 in zip(self.detections, self.frames_str):
    # ... (code for tracking objects across frames)
  
```

Filter Objects

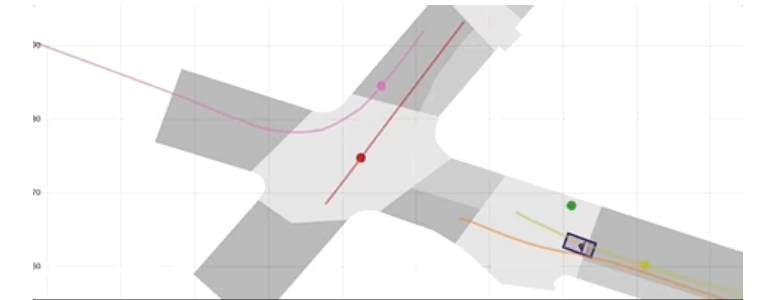
```

select c0.frameNum,
       c0.cameraId,
       c0.filename,
       t0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
       SegmentPolygon.polygon,
       valueAtTimestamp(
         t0.translations,
         c0.timestamp
       )
)
  
```

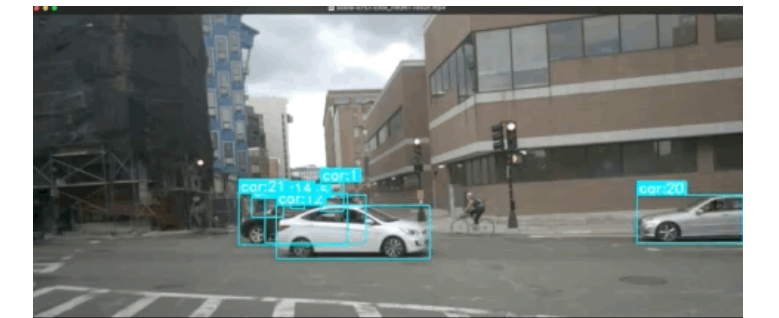
Format Outputs

FFmpeg
OpenCV

Analyze / Visualize

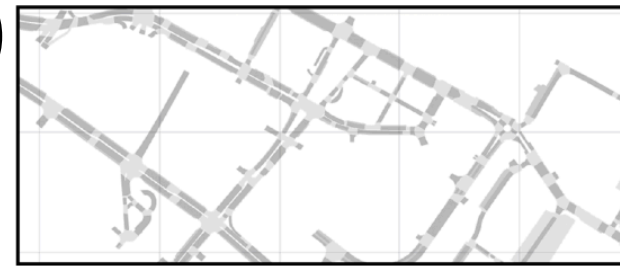


Annotate / Watch

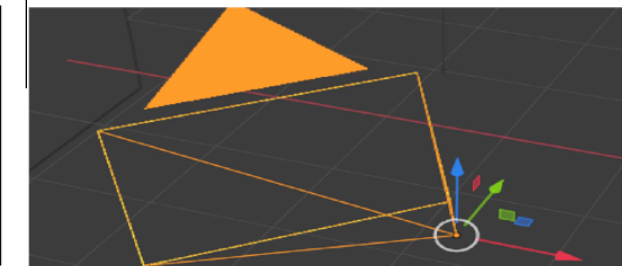


Typical Workflow (Pro)

Road Network



Camera Location / Rotation



Object Detector

```

with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(
        self.model.eval())

    assert isinstance(self.imgsz, List), type(self.
self.model.warmup(imgsz=(1, 3, *self.imgsz))) #

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imgsz, stride
im = im.transpose((2, 0, 1))[::-1] # HWC to
im = np.ascontiguousarray(im) # contiguous

        # t1 = time_sync()
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        # Inference
        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(
            pred,
            self.conf_thres,
            self.iou_thres,
            self.classes,
            self.agnostic_nms,
            max_det=self.max_det,
        )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(
det[:, :4] = scale_boxes(im.shape[2:], det[
        yield Detection2D(
            det, names, [DetectionFrame_idx, ord
  
```

3D Location Estimator

```

class LocationEstimator:
    def __init__(self, road_network, camera_location,
self.road_network = road_network
self.camera_location = camera_location

    def estimate_location(self, det):
        # Get bounding box coordinates
        x1, y1, x2, y2 = det[:, :4].cpu().numpy()

        # Find the closest road segment
        closest_road = self.road_network.closest_road(
x1, y1, x2, y2)

        # Estimate 3D location based on road geometry
        x3d, y3d, z3d = self.estimate_3d_location(
closest_road, x1, y1, x2, y2)

        return x3d, y3d, z3d
  
```

Object Tracker

```

class ObjectTracker:
    def __init__(self, detector, estimator, road_network):
self.detector = detector
self.estimator = estimator
self.road_network = road_network

    def track(self, frame_idx, im):
        # Detect objects
        det = self.detector.detect(im)

        # Estimate 3D location
        x3d, y3d, z3d = self.estimator.estimate_location(
det)

        # Track objects across frames
        track_id = self.track_id(det)

        return track_id, x3d, y3d, z3d
  
```

Filter Objects

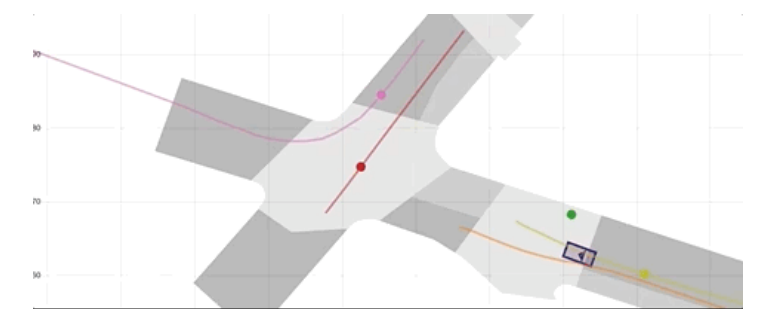
```

select c0.frameNum,
c0.cameraId,
c0.filename,
c0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
SegmentPolygon.polygon,
valueAtTimestamp(
t0.translations,
c0.timestamp
)
)
  
```

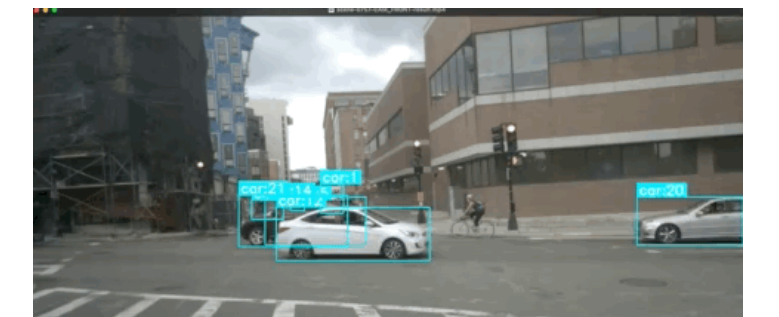
Format Outputs

FFmpeg
OpenCV

Analyze / Visualize



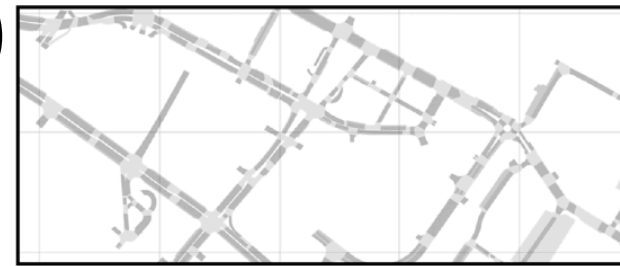
Annotate / Watch



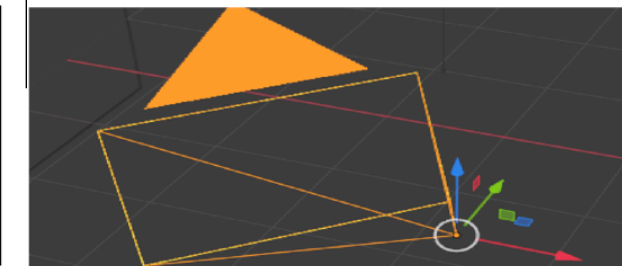
Error-Prone

Typical Workflow (Pro)

Road Network



Camera Location / Rotation



Object Detector

```

with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(
        self.model.eval())

    assert isinstance(self.imsz, List), type(self.
    self.model.warmup(imsz=(1, 3, *self.imsz)) #

    for frame_idx, im0 in enumerate(self.frames.str
    if isinstance(im0, Skip):
        yield skip
        continue

    im, _, _ = letterbox(im0, self.imsz, stride
    im = im.transpose((2, 0, 1))[::-1] # HWC t
    im = np.ascontiguousarray(im) # contiguous

    # t1 = time_sync()
    im = torch.from_numpy(im).to(self.device)
    im = im.half() if self.half else im.float()
    im /= 255.0 # 0 - 255 to 0.0 - 1.0
    if len(im.shape) == 3:
        im = im[None] # expand for batch dim

    # Inference
    pred = self.model(im, augment=self.augment)

    # Apply NMS
    pred = non_max_suppression(
        pred,
        self.conf_thres,
        self.iou_thres,
        self.classes,
        self.agnostic_nms,
        max_det=self.max_det,
    )

    # Process detections
    assert isinstance(pred, List), type(pred)
    assert len(pred) == 1, len(pred)
    det = pred[0]
    assert isinstance(det, torch.Tensor), type
    det[:, :4] = scale_boxes(im.shape[2:], det[
    det, names, [DetectionFrame, ...], ord
  
```

3D Location Estimator

```

class LocationEstimator:
    def __init__(self, road_network, camera_config):
        self.road_network = road_network
        self.camera_config = camera_config

    def estimate_location(self, detections):
        # ... (omitted code) ...
  
```

Object Tracker

```

saved_detections: list[list, torch.Tensor] = []
class ListList: None = None
FFmpeg = None

for detection, im0 in zip(
    self.detections, self.frames.strided):
    structure,
    structure,
    )
    if not isinstance(detection, Skip):
        assert not isinstance(im0, Skip), type(im0)
        im0 = im0.copy()

    # ... (omitted code) ...
  
```

Filter Objects

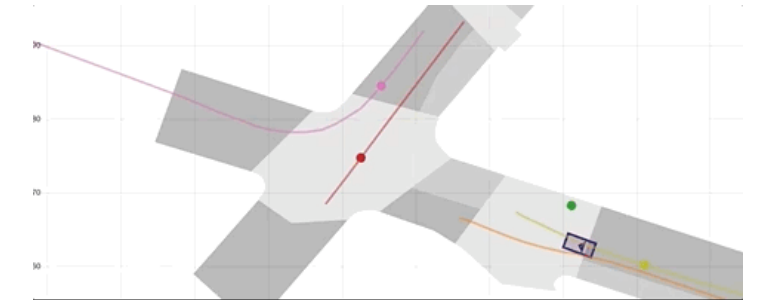
```

select c0.frameNum,
       c0.cameraId,
       c0.filename,
       c0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
    SegmentPolygon.polygon,
    valueAtTimestamp(
        t0.translations,
        c0.timestamp
    )
)
  
```

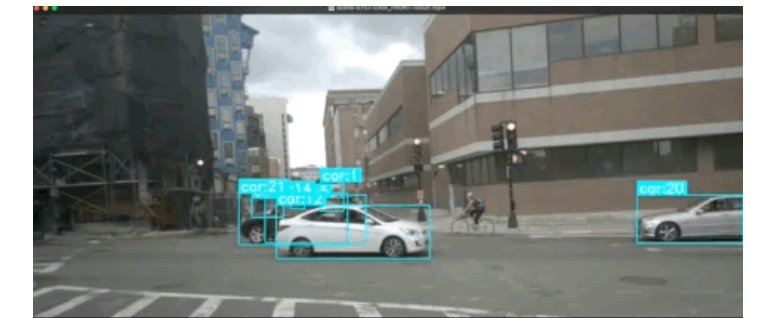
Format Outputs

FFmpeg
OpenCV

Analyze / Visualize



Annotate / Watch

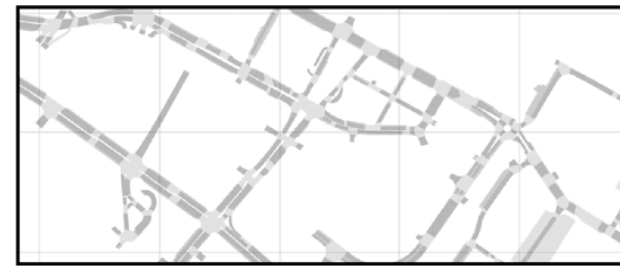


Expensive ⌚ + 💰

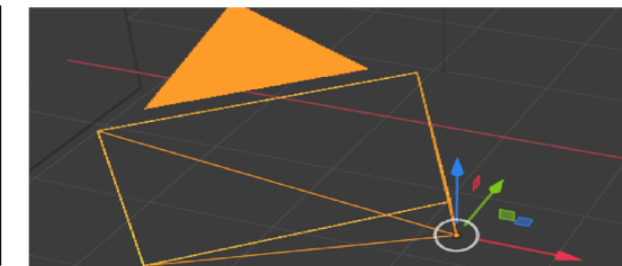
Error-Prone

Spatialize Workflow

Road Network



Camera Location / Rotation



Object Detector

```

with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(self.model.eval())

    assert isinstance(self.imsz, List), type(self.imsz)
    self.model.warmup(imsz=(1, 3, *self.imsz)) #

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imsz, stride=self.stride)
        im = im.transpose(2, 0, 1) # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det, )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4], im.shape[2:], )
        yield Detection2D(det, names, [DetectionFrame, ], ord=)
  
```

3D Location Estimator

```

class LocationEstimator:
    def __init__(self, model, camera_config):
        self.model = model
        self.camera_config = camera_config

    def estimate_location(self, frame_idx, im0):
        # Preprocess image
        im, _, _ = letterbox(im0, self.imsz, stride=self.stride)
        im = im.transpose(2, 0, 1)
        im = np.ascontiguousarray(im)

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None]

        pred = self.model(im)

        # Process detections
        assert len(pred) == 1
        det = pred[0]

        # Estimate location
        # ... (detailed code for location estimation) ...

        return location
  
```

Object Tracker

```

class ObjectTracker:
    def __init__(self, model, camera_config):
        self.model = model
        self.camera_config = camera_config

    def track_objects(self, frame_idx, im0):
        # Detect objects
        det = self.location_estimator.estimate_location(frame_idx, im0)

        # Track objects
        # ... (detailed code for object tracking) ...

        return tracks
  
```

Filter Objects

```

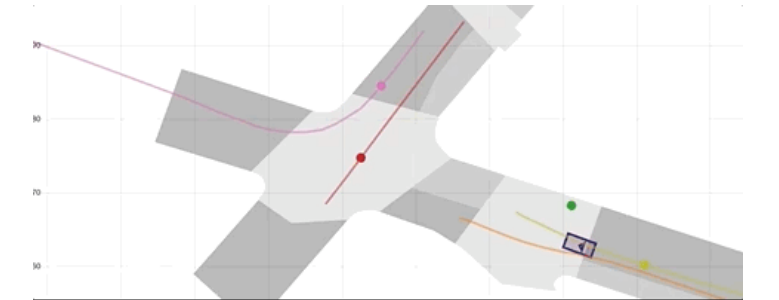
select c0.frameNum,
       c0.cameraId,
       c0.filename,
       t0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
       SegmentPolygon.polygon,
       valueAtTimestamp(
           t0.translations,
           c0.timestamp
       )
)
  
```

Format Outputs

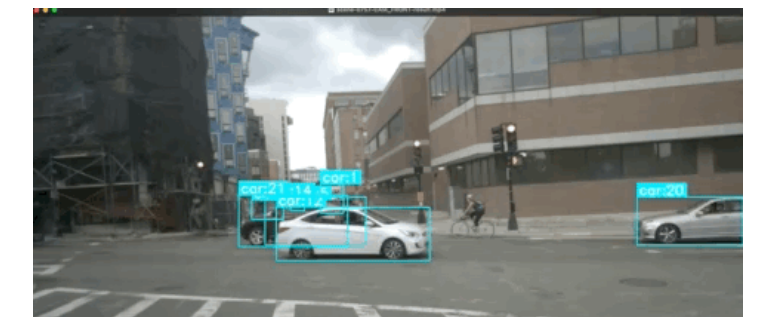
FFmpeg

OpenCV

Analyze / Visualize



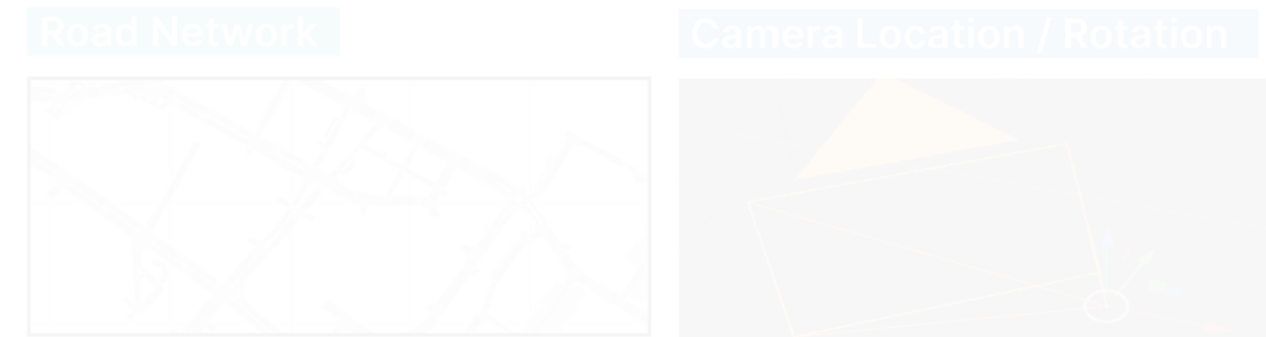
Annotate / Watch



Error-Prone

Expensive ⌚ + 💰

Spatialyze Workflow



1 The users only need to describe “what” their videos of interest look like instead of “how” to get them.

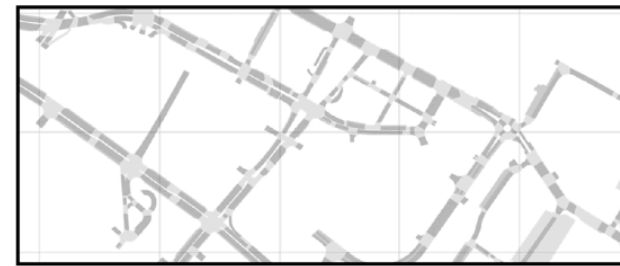
2 Spatialyze leverages users’ description of their videos of interest and the geospatial metadata to optimize for faster video processing.

Expensive ⌚ + 💰

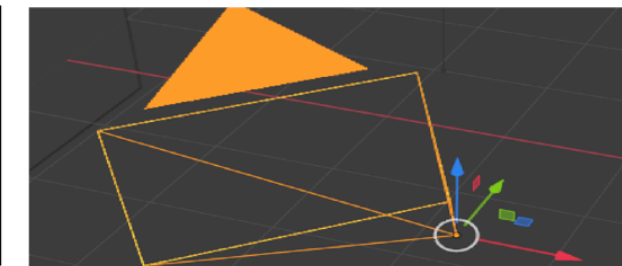
Error-Prone

Spatialize Workflow

Road Network



Camera Location / Rotation



Object Detector

```

with torch.no_grad():
    _names = self.model.names

    assert isinstance(_names, dict), type(_names)
    names: list[str] = class_mapping_to_list_names(self.model.eval())

    assert isinstance(self.imsz, List), type(self.imsz)
    self.model.warmup(imsz=(1, 3, *self.imsz)) #

    for frame_idx, im0 in enumerate(self.frames_str):
        if isinstance(im0, Skip):
            yield skip
            continue

        im, _, _ = letterbox(im0, self.imsz, stride=self.stride)
        im = im.transpose(2, 0, 1)[::-1] # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det, )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4], im.shape[2:], )
        yield Detection2D(det, names, [DetectionFrame, ], ord=)
  
```

3D Location Estimator

```

class LocationEstimator:
    def __init__(self, model, device, camera_config):
        self.model = model
        self.device = device
        self.camera_config = camera_config

    def estimate_location(self, frame_idx, im0, camera_config):
        # Preprocess image
        im, _, _ = letterbox(im0, self.imsz, stride=self.stride)
        im = im.transpose(2, 0, 1)[::-1] # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det, )

        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4], im.shape[2:], )
        yield Detection2D(det, names, [DetectionFrame, ], ord=)
  
```

Object Tracker

```

saved_detections: list[list, torch.Tensor] = []
class Tracker:
    def __init__(self, model, device, camera_config):
        self.model = model
        self.device = device
        self.camera_config = camera_config

    def track(self, frame_idx, im0, camera_config):
        # Preprocess image
        im, _, _ = letterbox(im0, self.imsz, stride=self.stride)
        im = im.transpose(2, 0, 1)[::-1] # HWC to CHW
        im = np.ascontiguousarray(im) # contiguous

        # Inference
        im = torch.from_numpy(im).to(self.device)
        im = im.half() if self.half else im.float()
        if len(im.shape) == 3:
            im = im[None] # expand for batch dim

        pred = self.model(im, augment=self.augment)

        # Apply NMS
        pred = non_max_suppression(pred, self.conf_thres, self.iou_thres, self.classes, self.agnostic_nms, max_det=self.max_det, )



        # Process detections
        assert isinstance(pred, List), type(pred)
        assert len(pred) == 1, len(pred)
        det = pred[0]
        assert isinstance(det, torch.Tensor), type(det)
        det[:, :4] = scale_boxes(im.shape[2:], det[:, :4], im.shape[2:], )
        yield Detection2D(det, names, [DetectionFrame, ], ord=)
  
```

Filter Objects

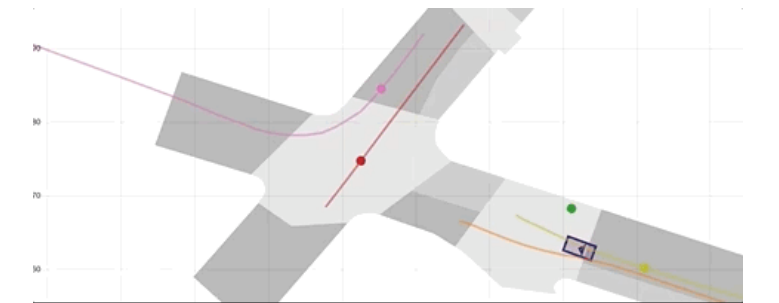
```

select c0.frameNum,
       c0.cameraId,
       c0.filename,
       t0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
       SegmentPolygon.polygon,
       valueAtTimestamp(
           t0.translations,
           c0.timestamp
       )
)
  
```

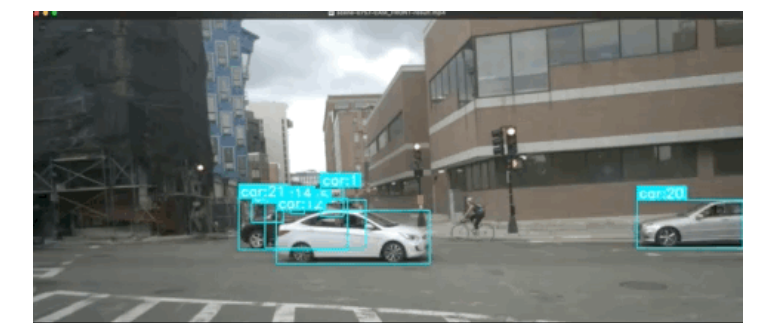
Format Outputs

Analyze / Visualize



Annotate / Watch

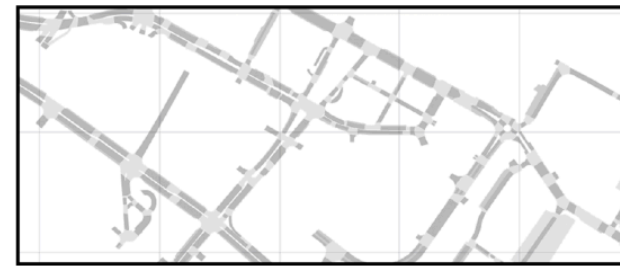


Expensive ⌚ + 💰

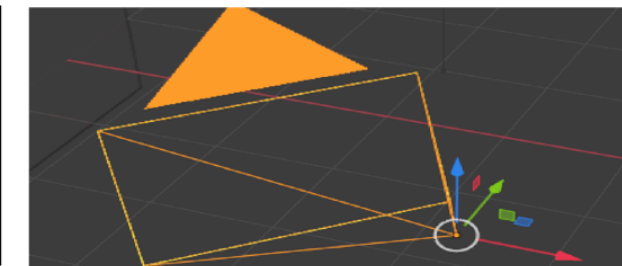
Error-Prone

Spatialize Workflow

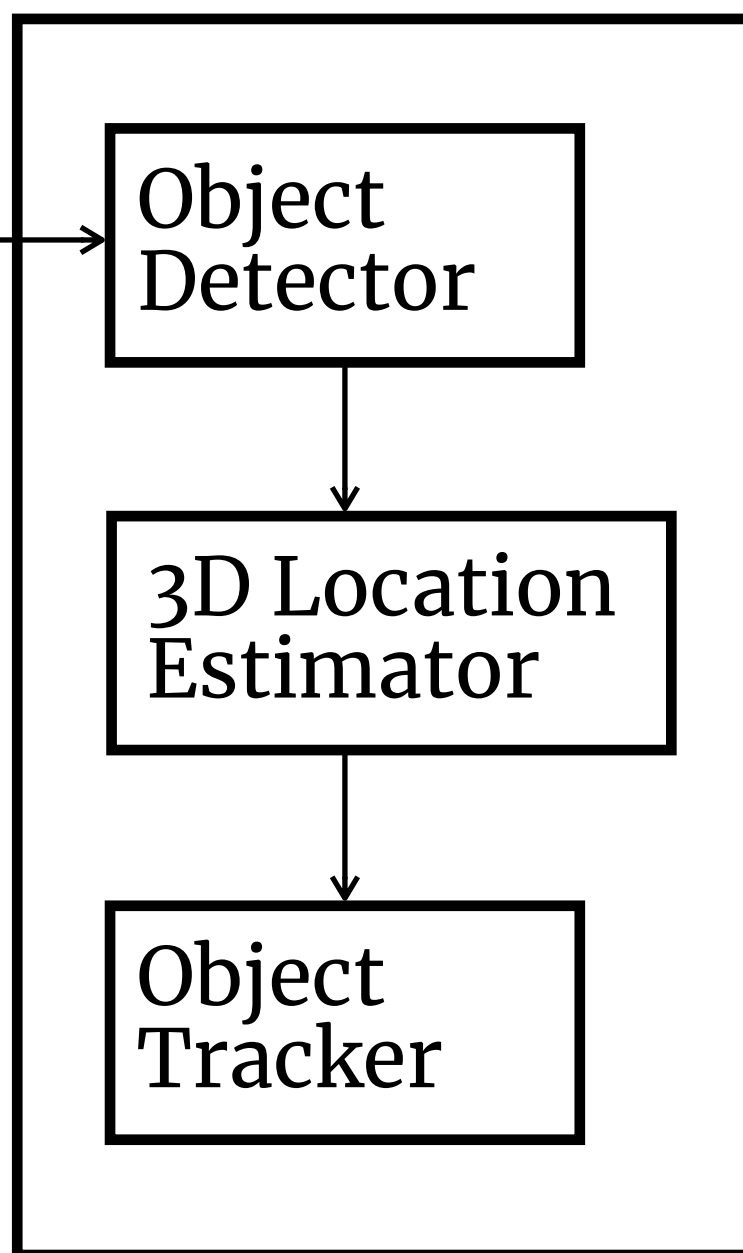
Road Network



Camera Location / Rotation



SPATIALYZE



Filter Objects

```
select c0.frameNum,
c0.cameraId,
c0.filename,
t0.itemid
FROM Camera as c0
JOIN Item_Trajectory as t0
ON c0.timestamp <@ t0.translations::period
AND c0.cameraId = t0.cameraId
JOIN SegmentPolygon
WHERE ST_Contains(
SegmentPolygon.polygon,
valueAtTimestamp(
t0.translations,
c0.timestamp
)
)
```

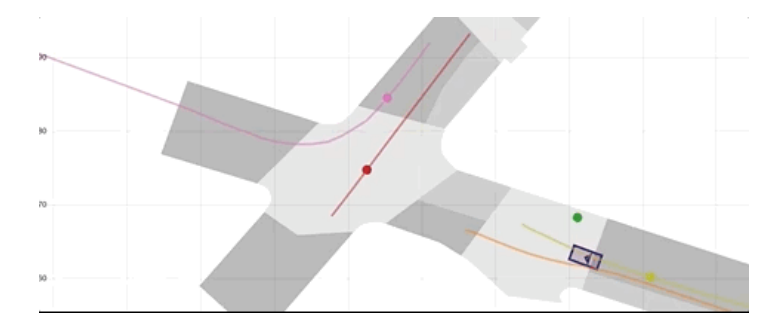


Format Outputs

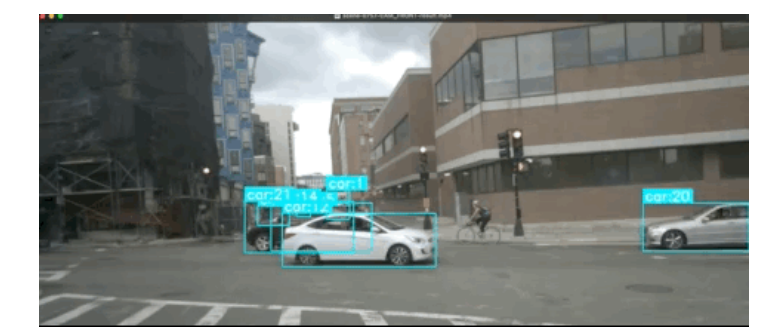
FFmpeg

OpenCV

Analyze / Visualize



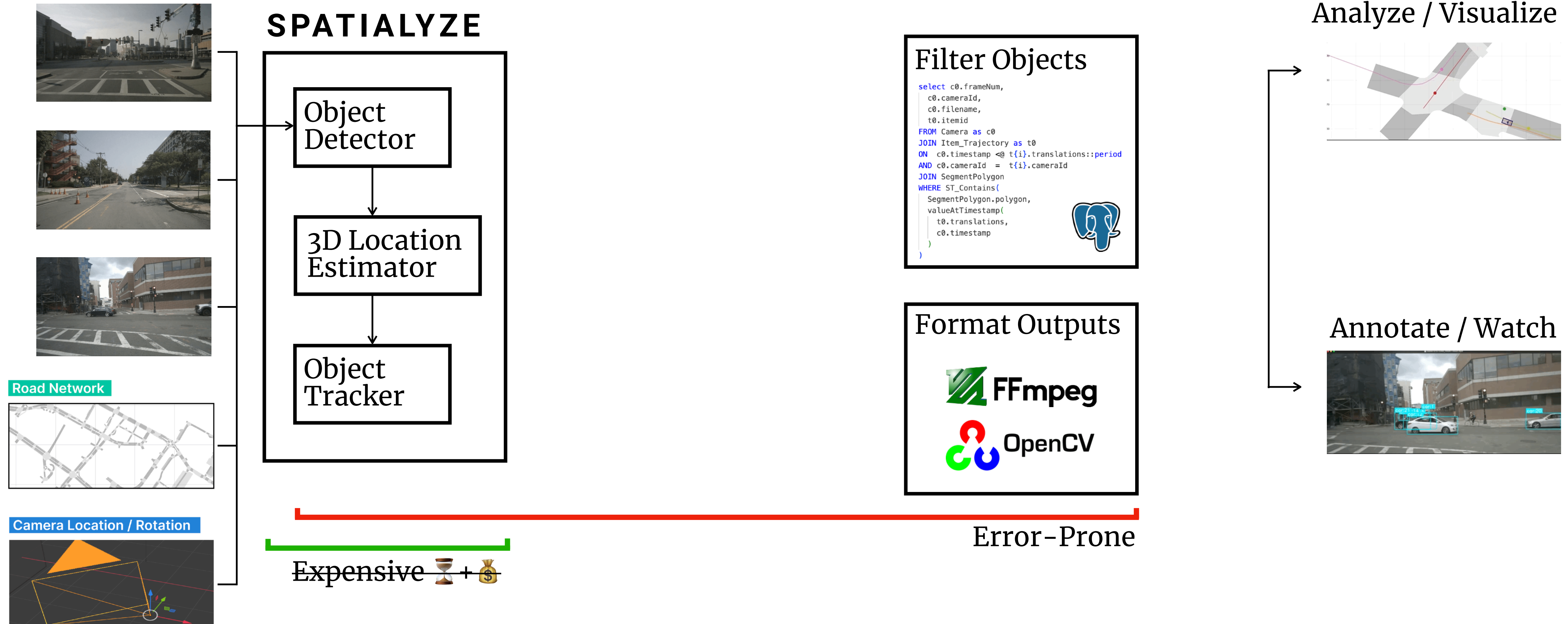
Annotate / Watch



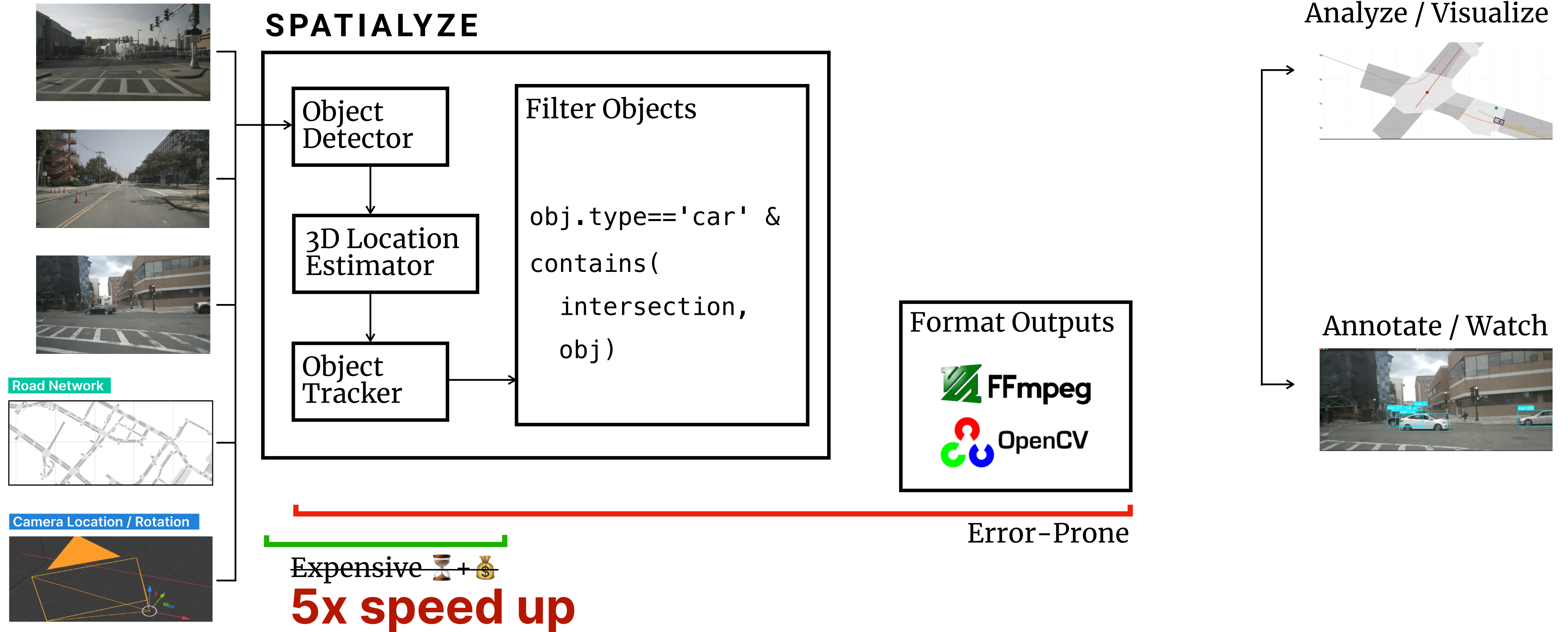
Expensive ⌚ + 💰

Error-Prone

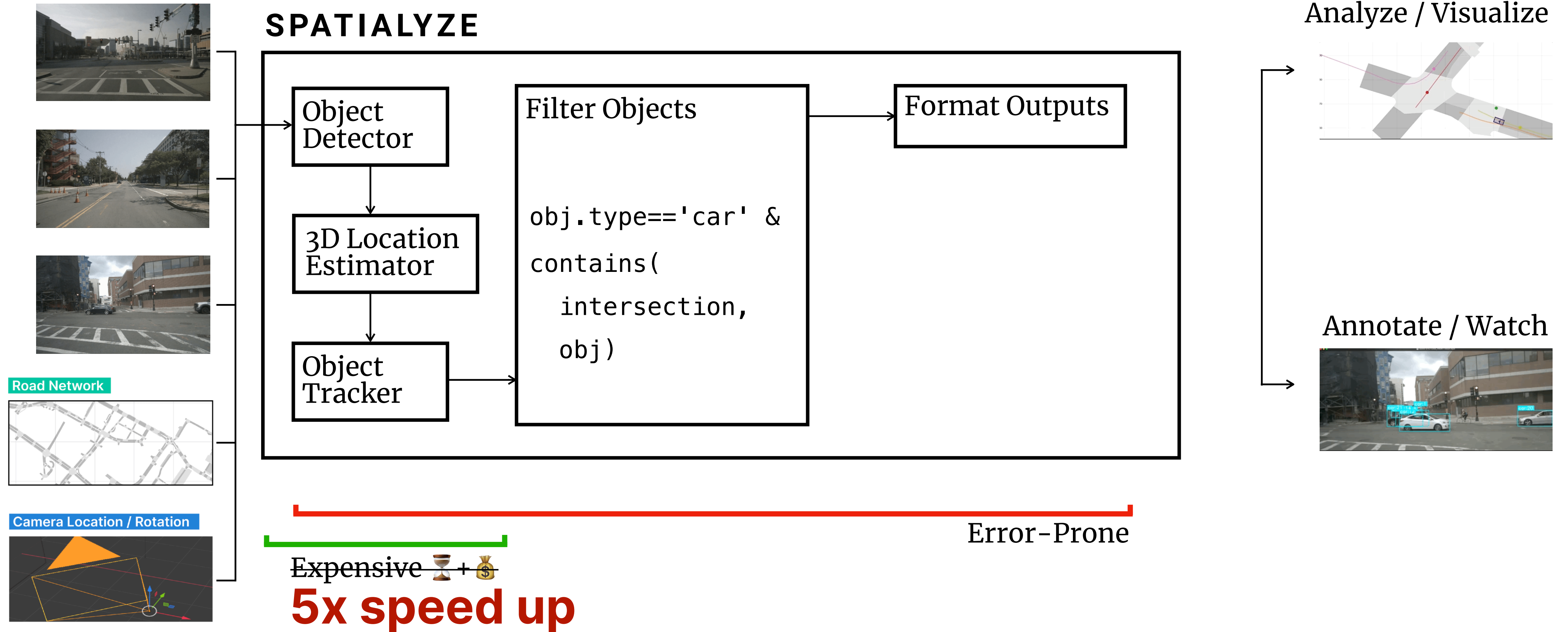
Spatialize Workflow



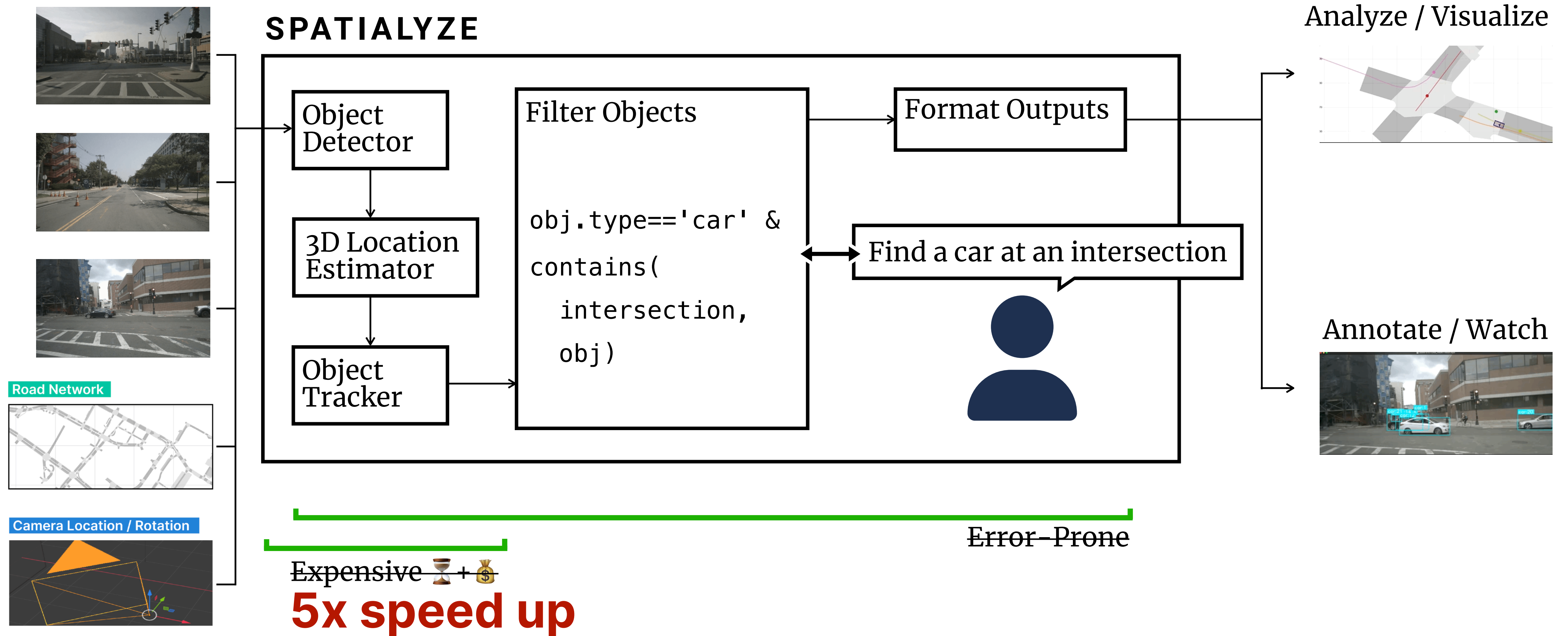
Spatialize Workflow



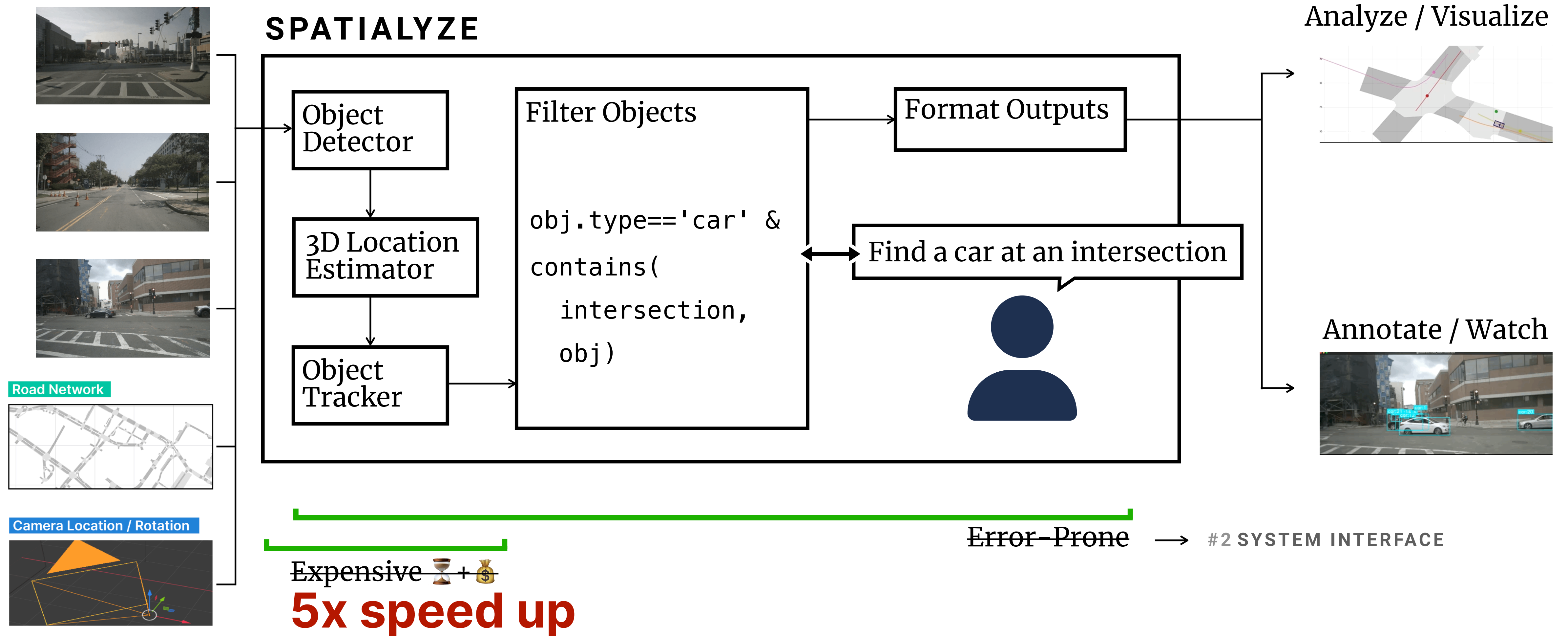
Spatialize Workflow



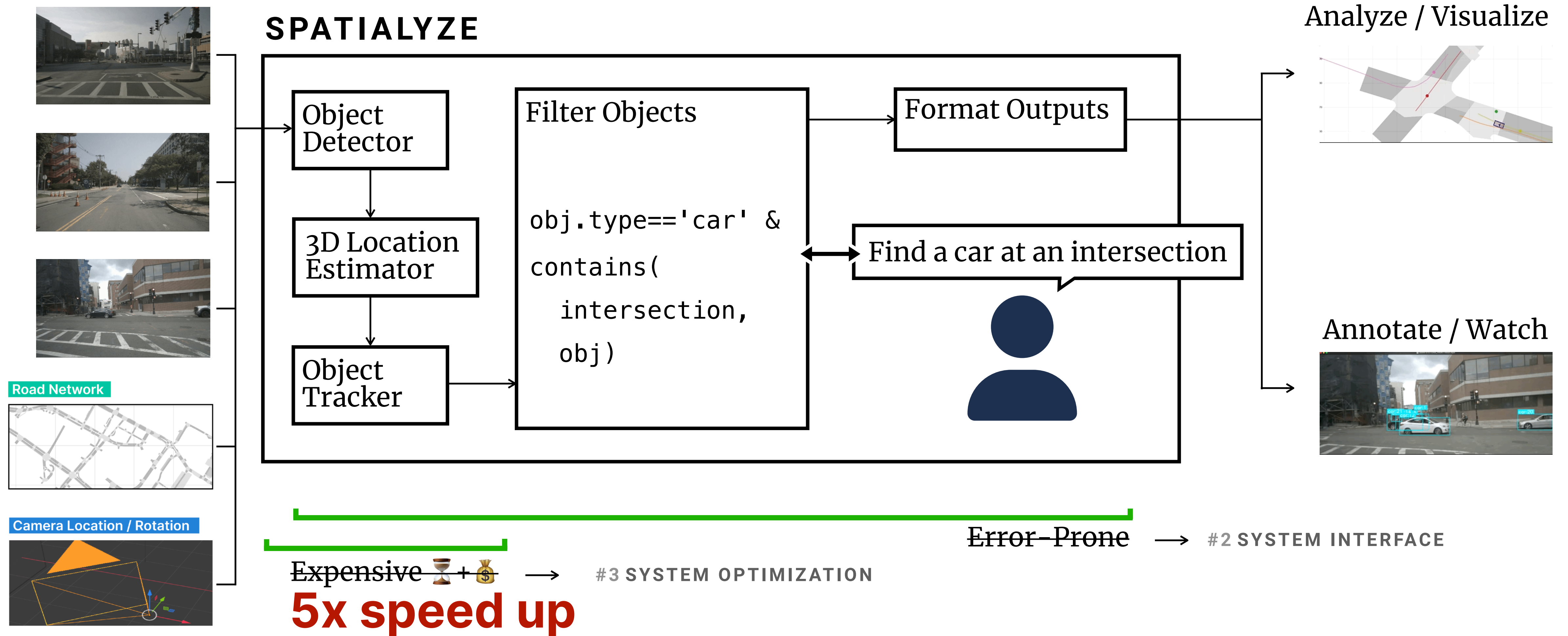
Spatialize Workflow



Spatialize Workflow



Spatialize Workflow



SPATIALYZE

#2 **System Interface**

- * Data Model
- * Programming Model

#2 SYSTEM INTERFACE

Data Model



Data Model

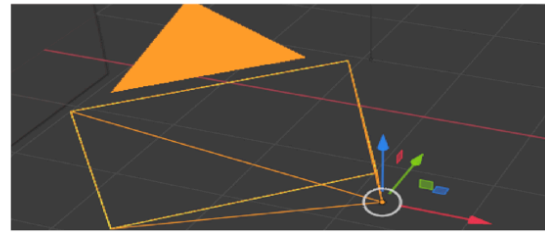


Data Model

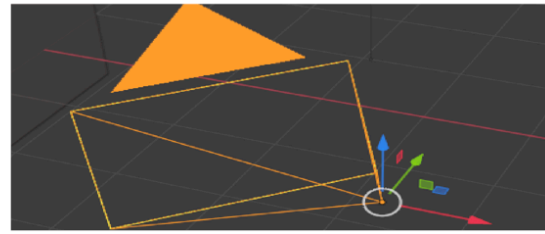


Data Model

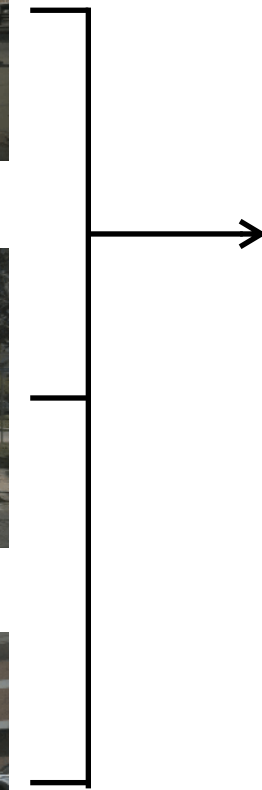
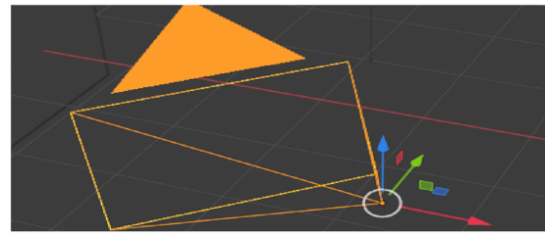
Camera Location / Rotation



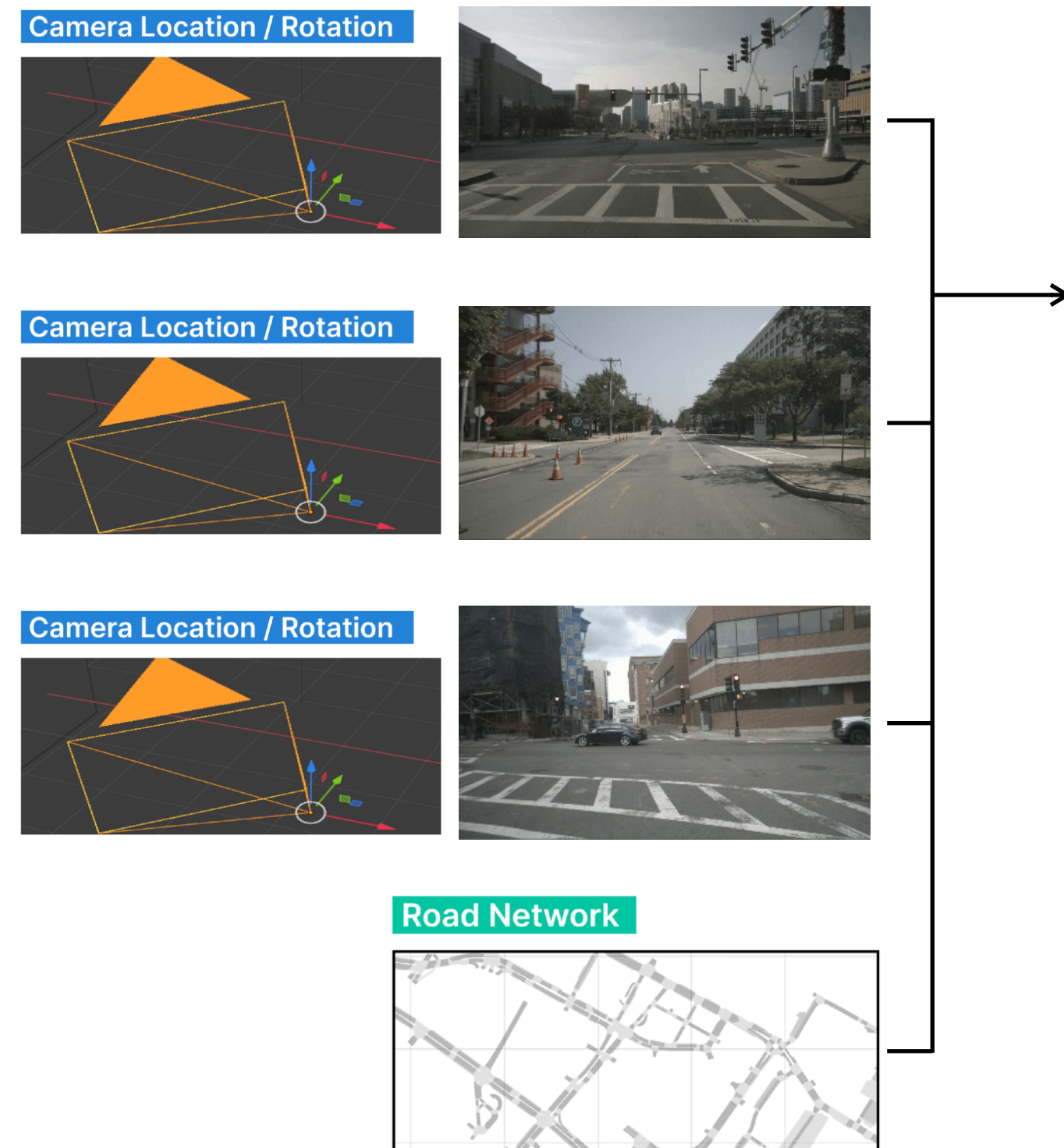
Camera Location / Rotation



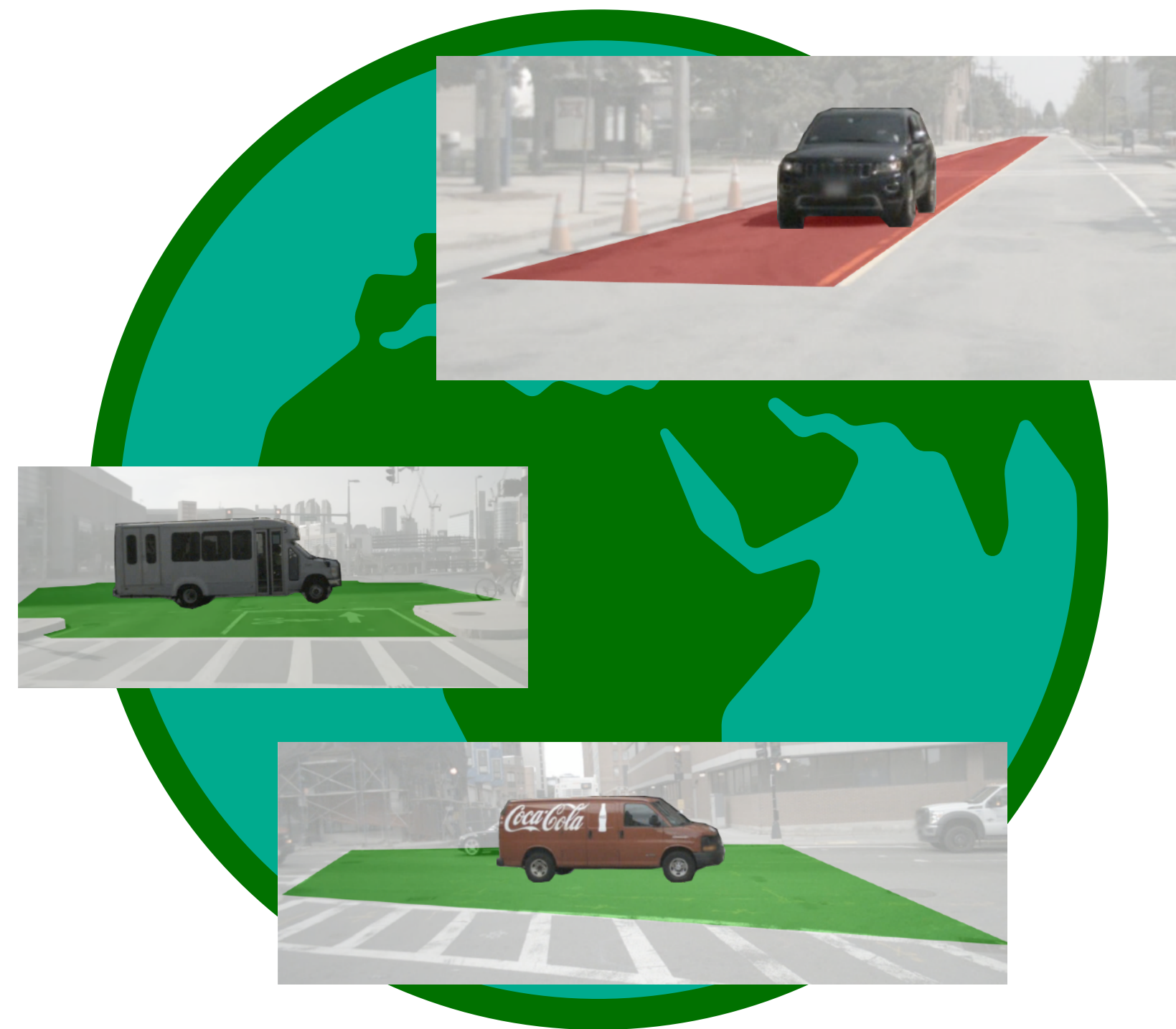
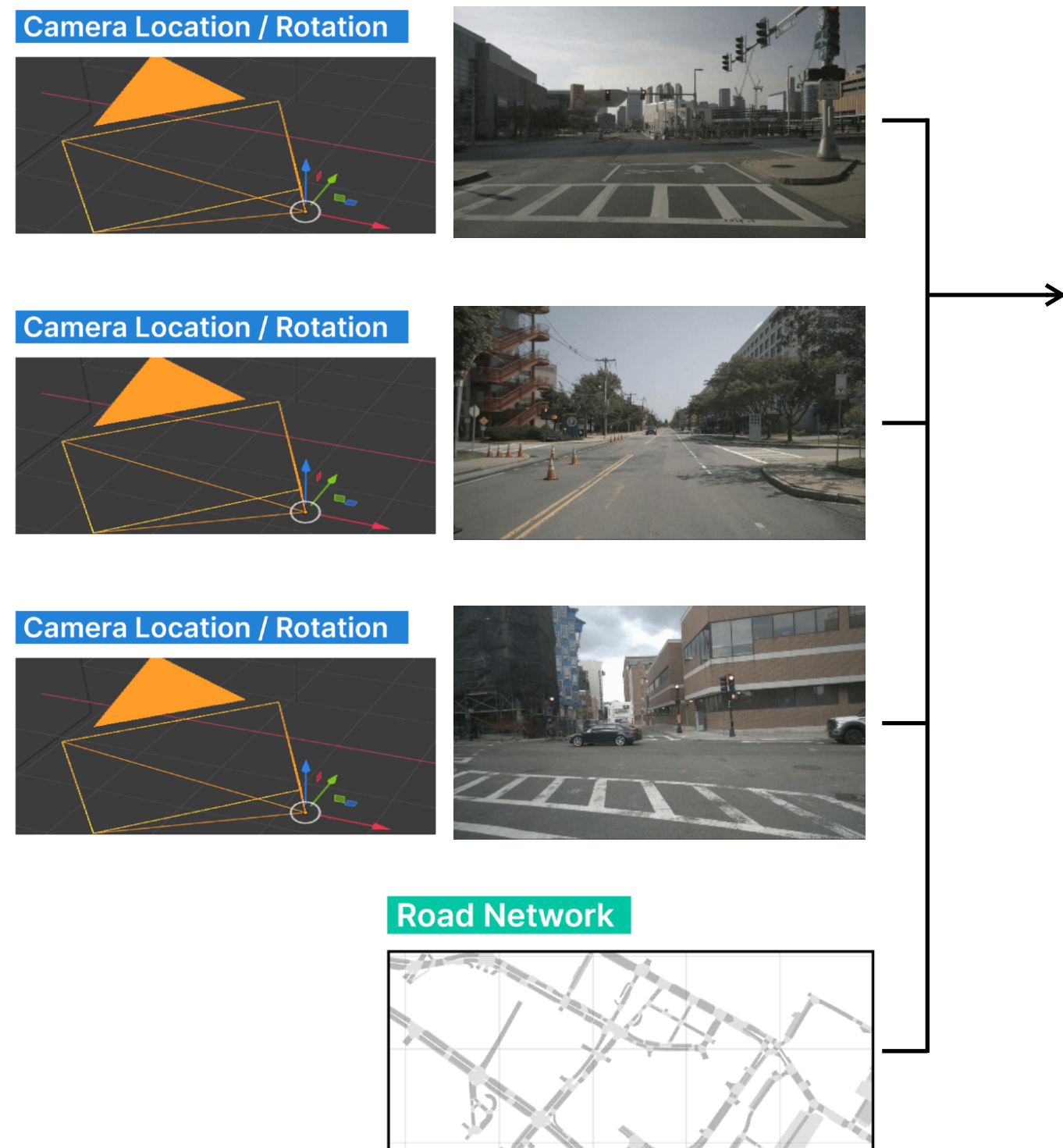
Camera Location / Rotation



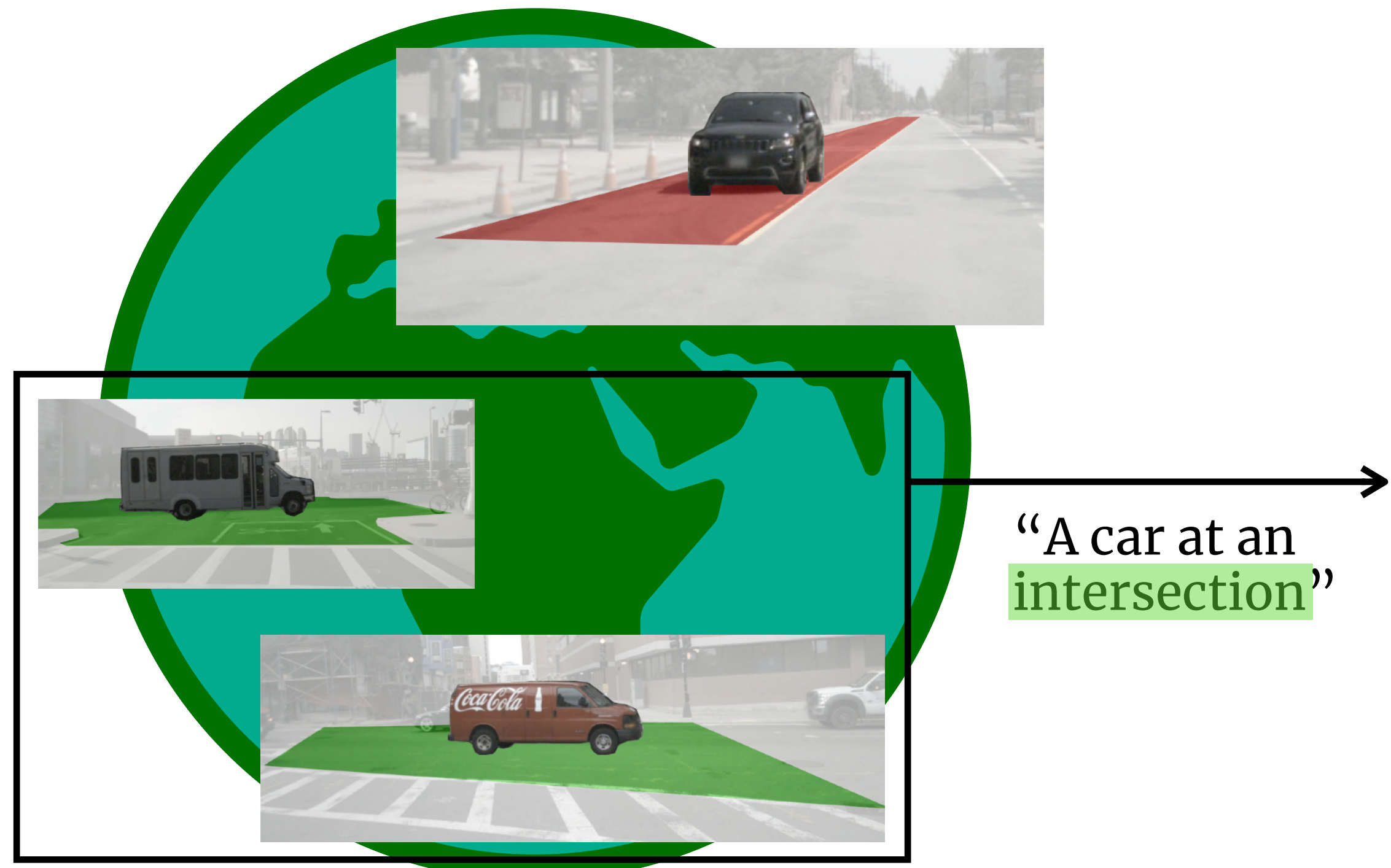
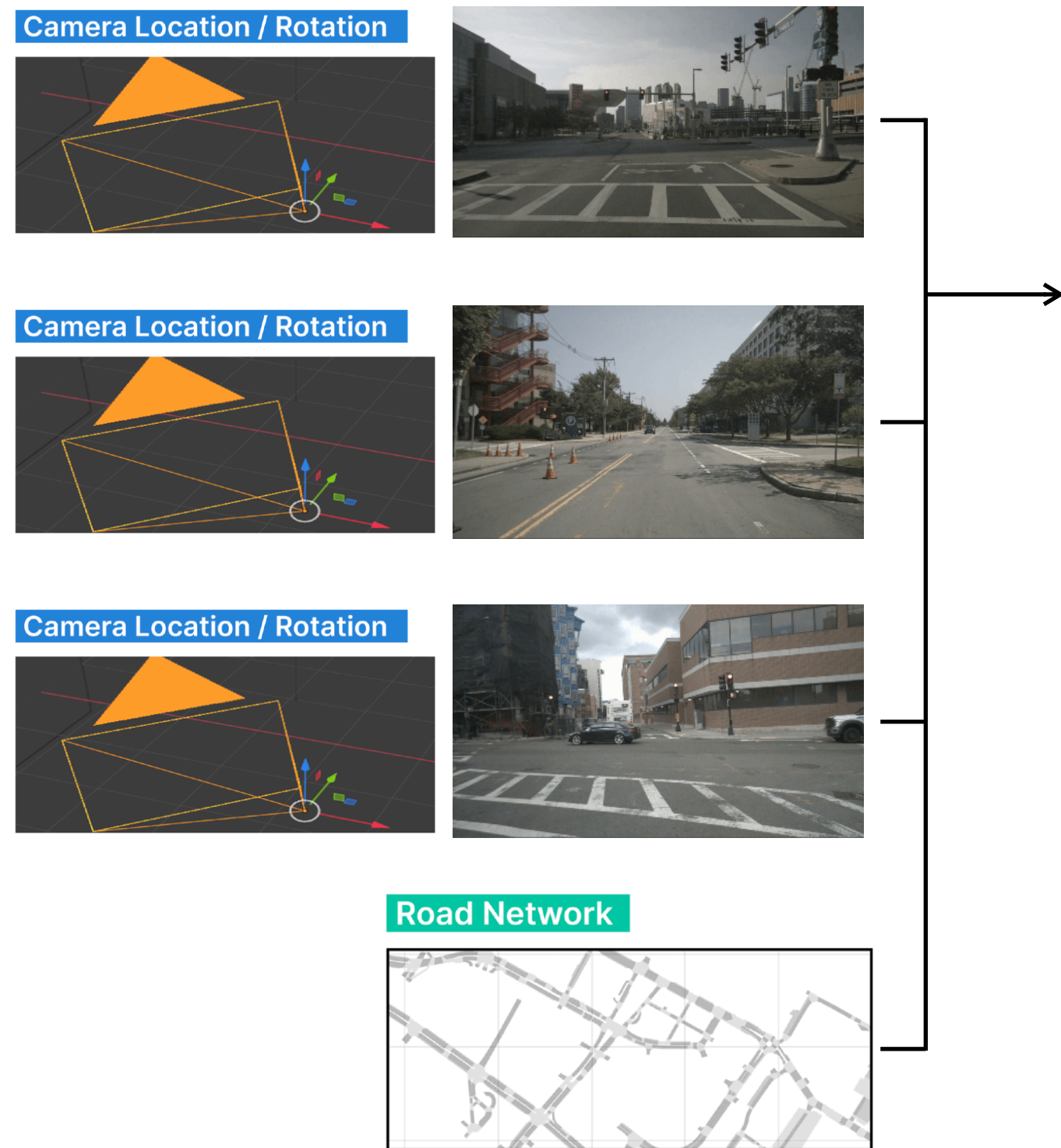
Data Model



Data Model

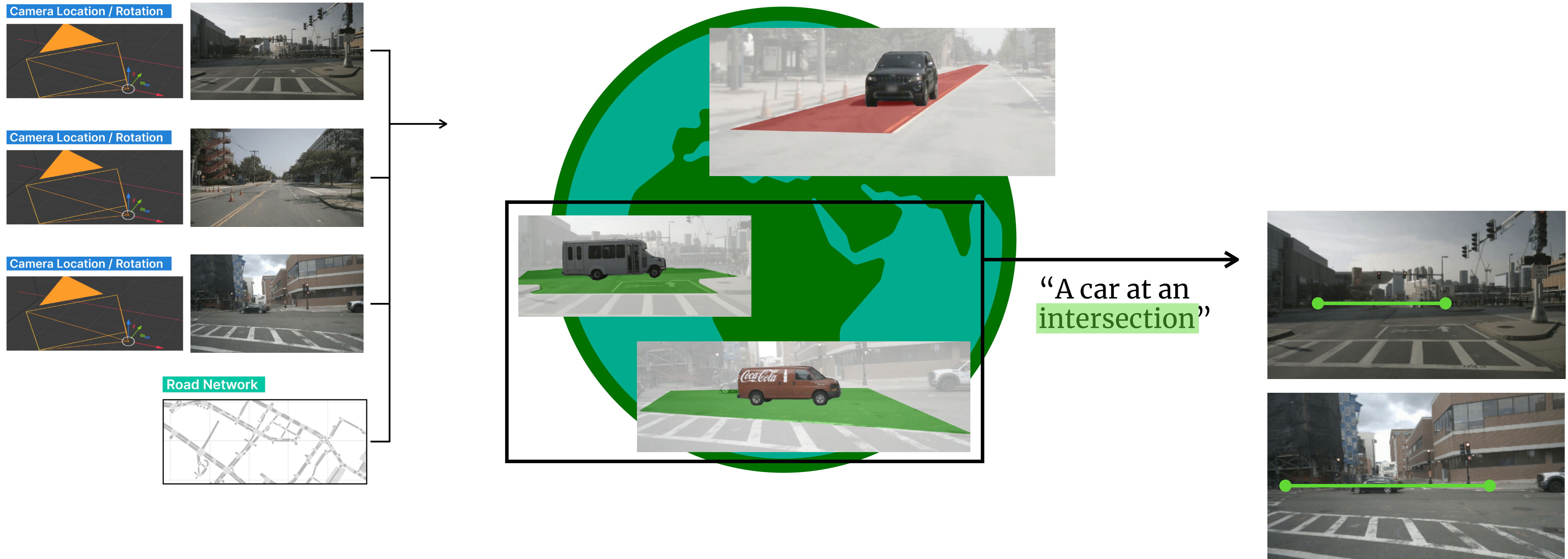


Data Model

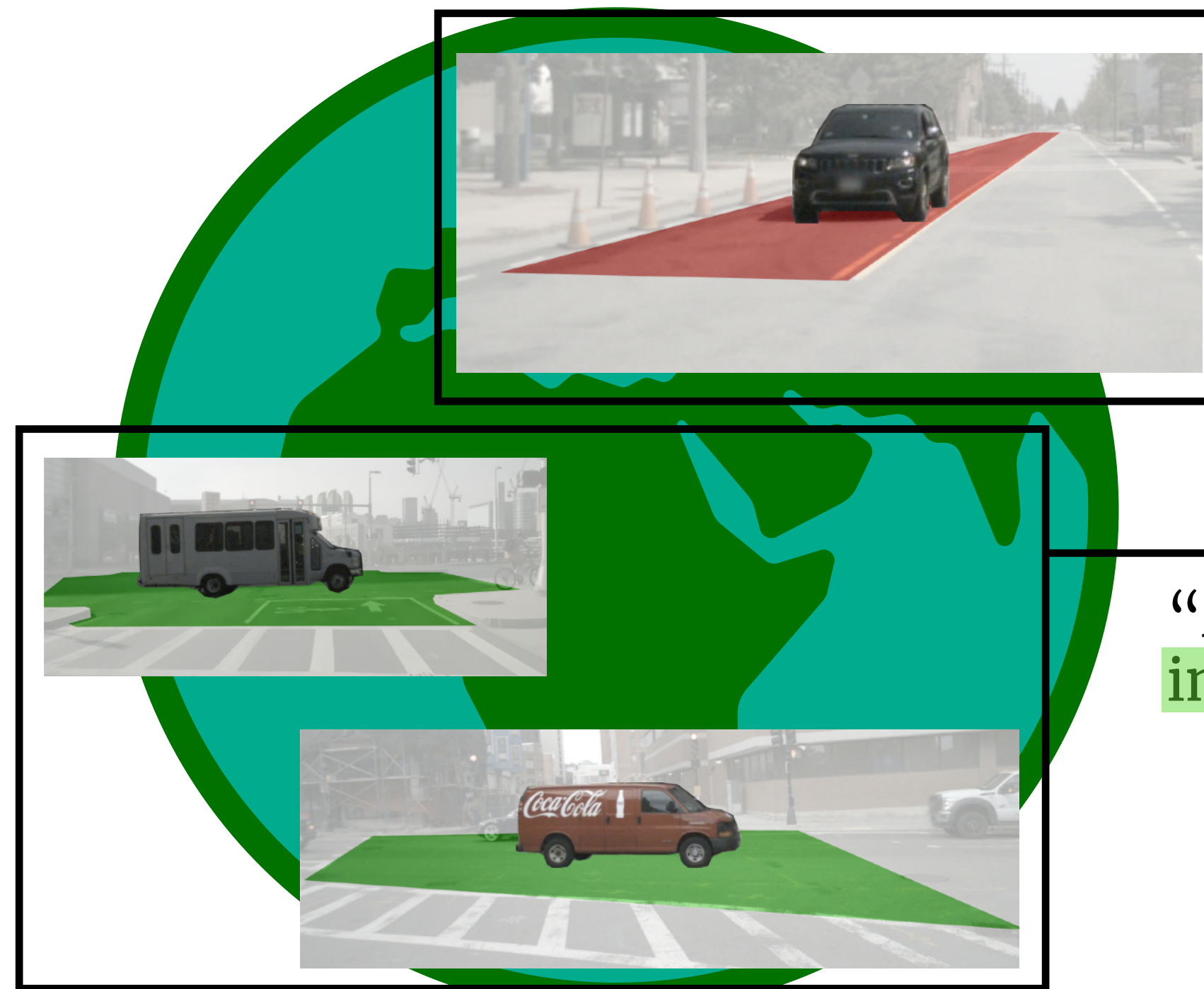
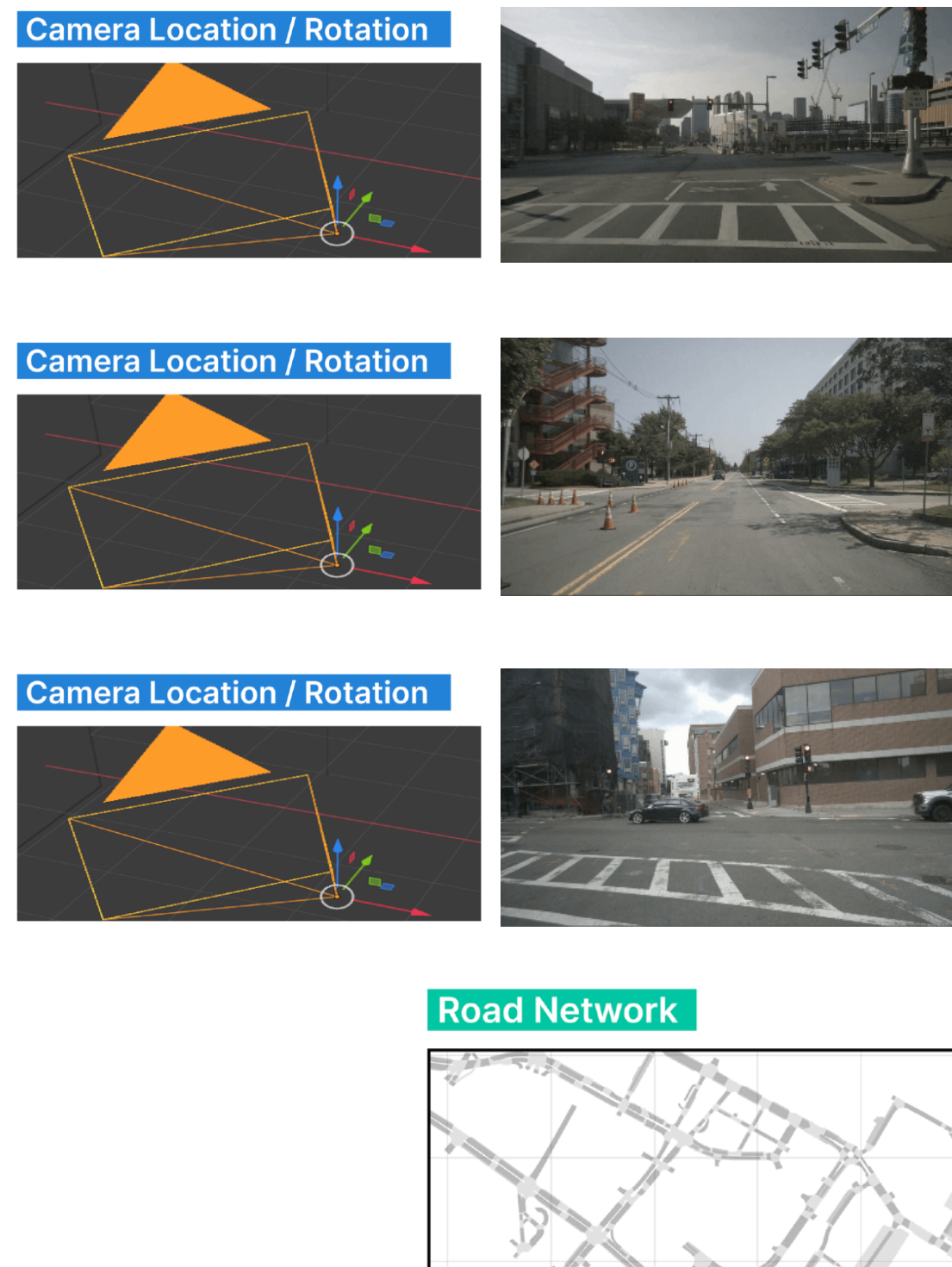


"A car at an intersection"

Data Model

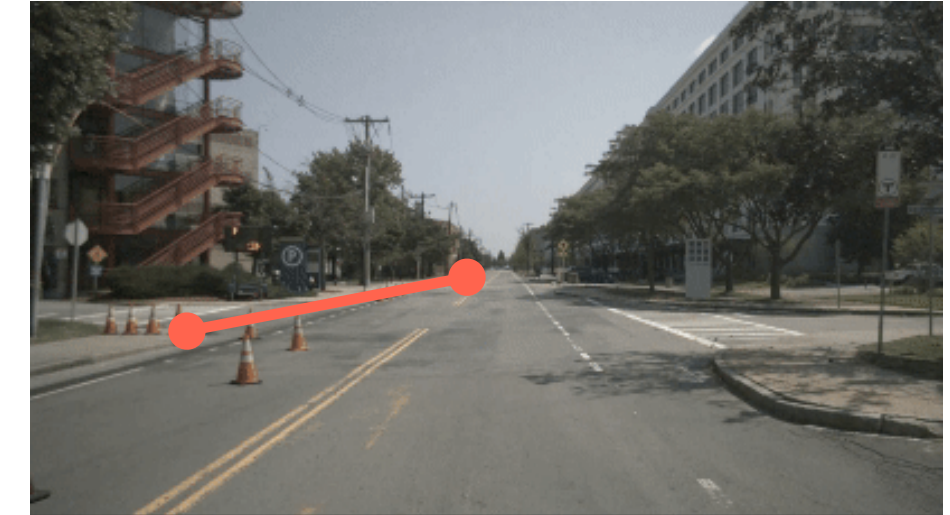


Data Model



“A car on a lane”

“A car at an intersection”



Programming Model

Build → *Filter* → *Observe*

Programming Model



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

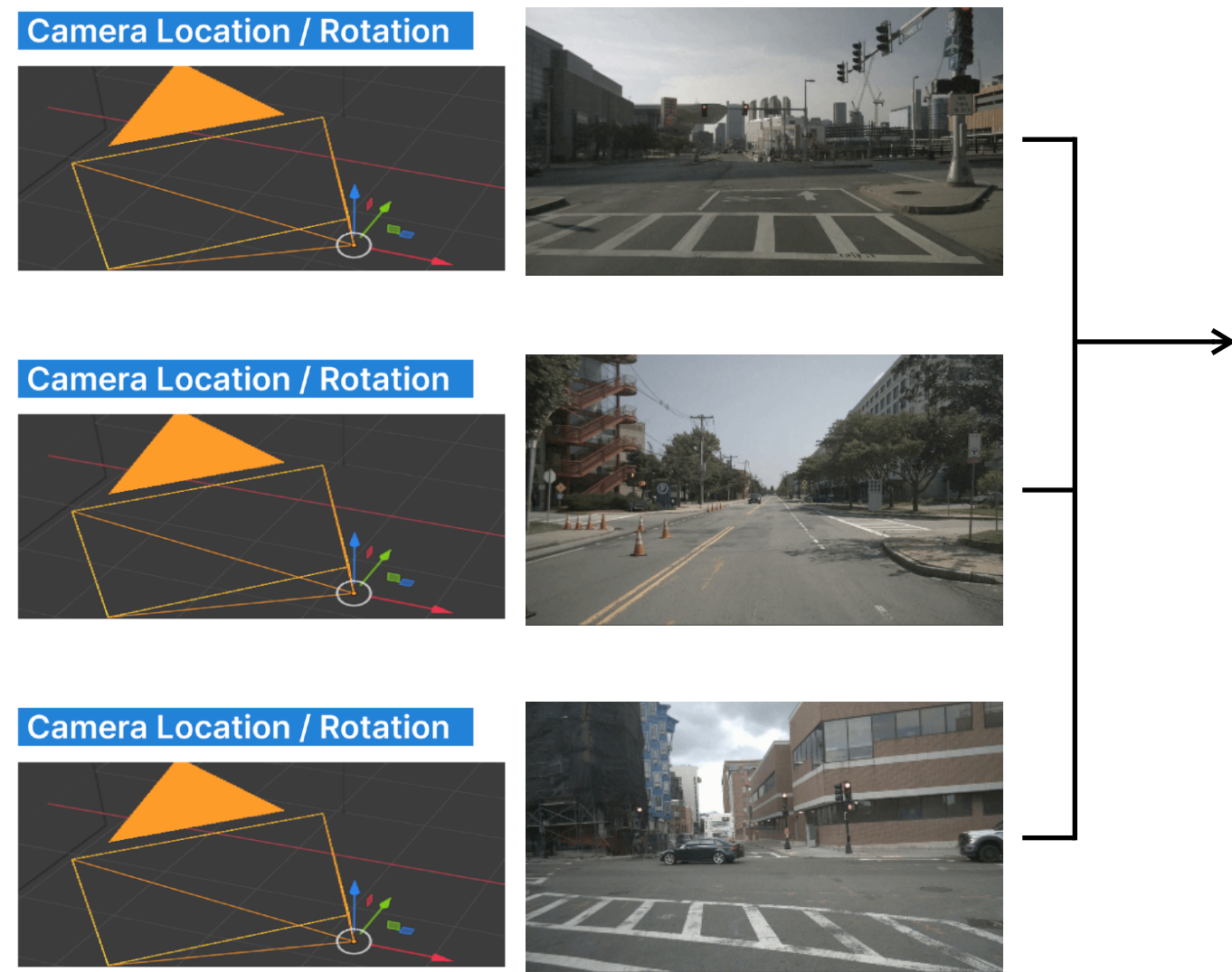
Filter

```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
        contains(i, o))
```

Observe

```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Programming Model



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

Filter

```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
        contains(i, o))
```

Observe

```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Programming Model



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

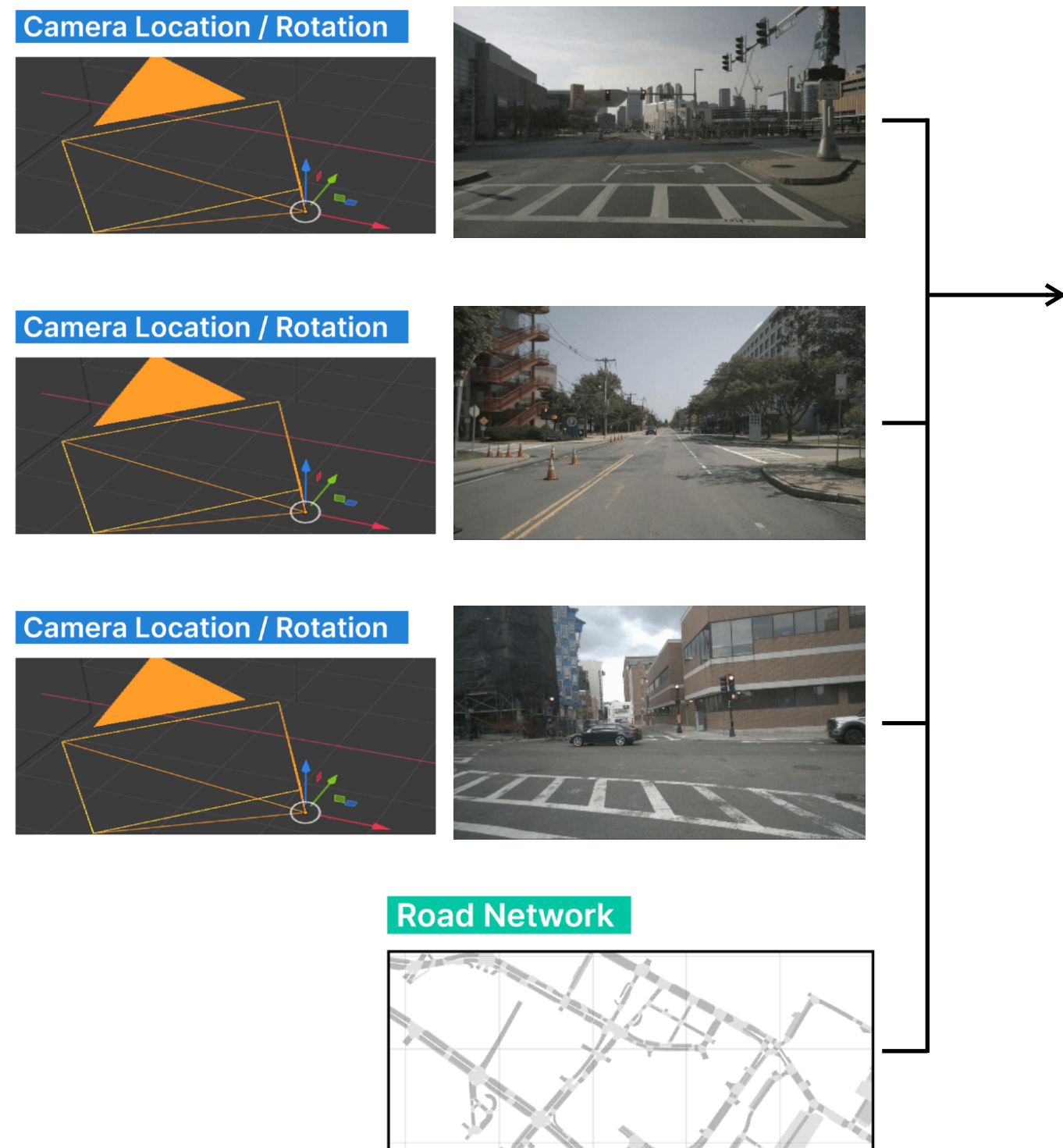
Filter

```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
         contains(i, o))
```

Observe

```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Programming Model



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

Filter

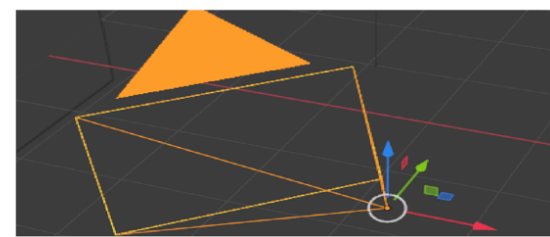
```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
         contains(i, o))
```

Observe

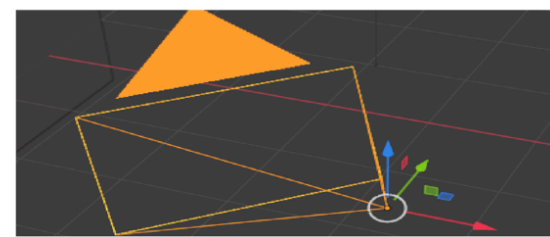
```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Programming Model

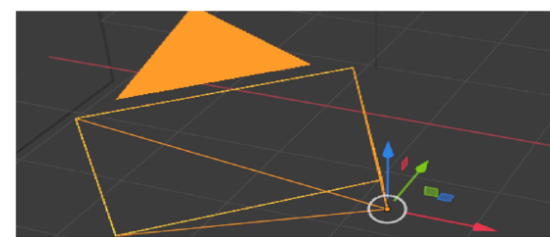
Camera Location / Rotation



Camera Location / Rotation



Camera Location / Rotation



Road Network



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

Filter

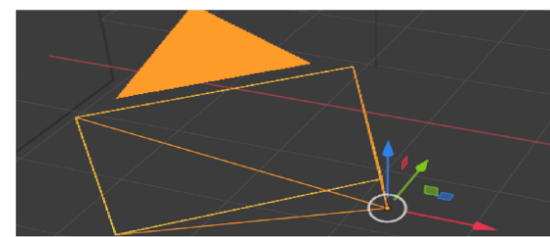
```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
         contains(i, o))
```

Observe

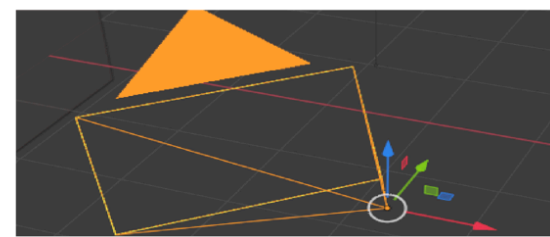
```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Programming Model

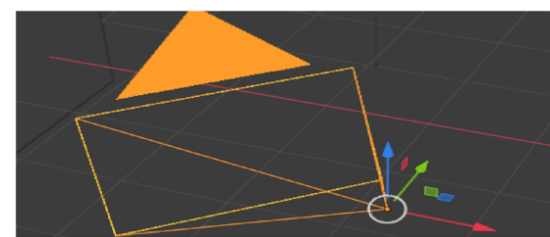
Camera Location / Rotation



Camera Location / Rotation



Camera Location / Rotation



Road Network



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

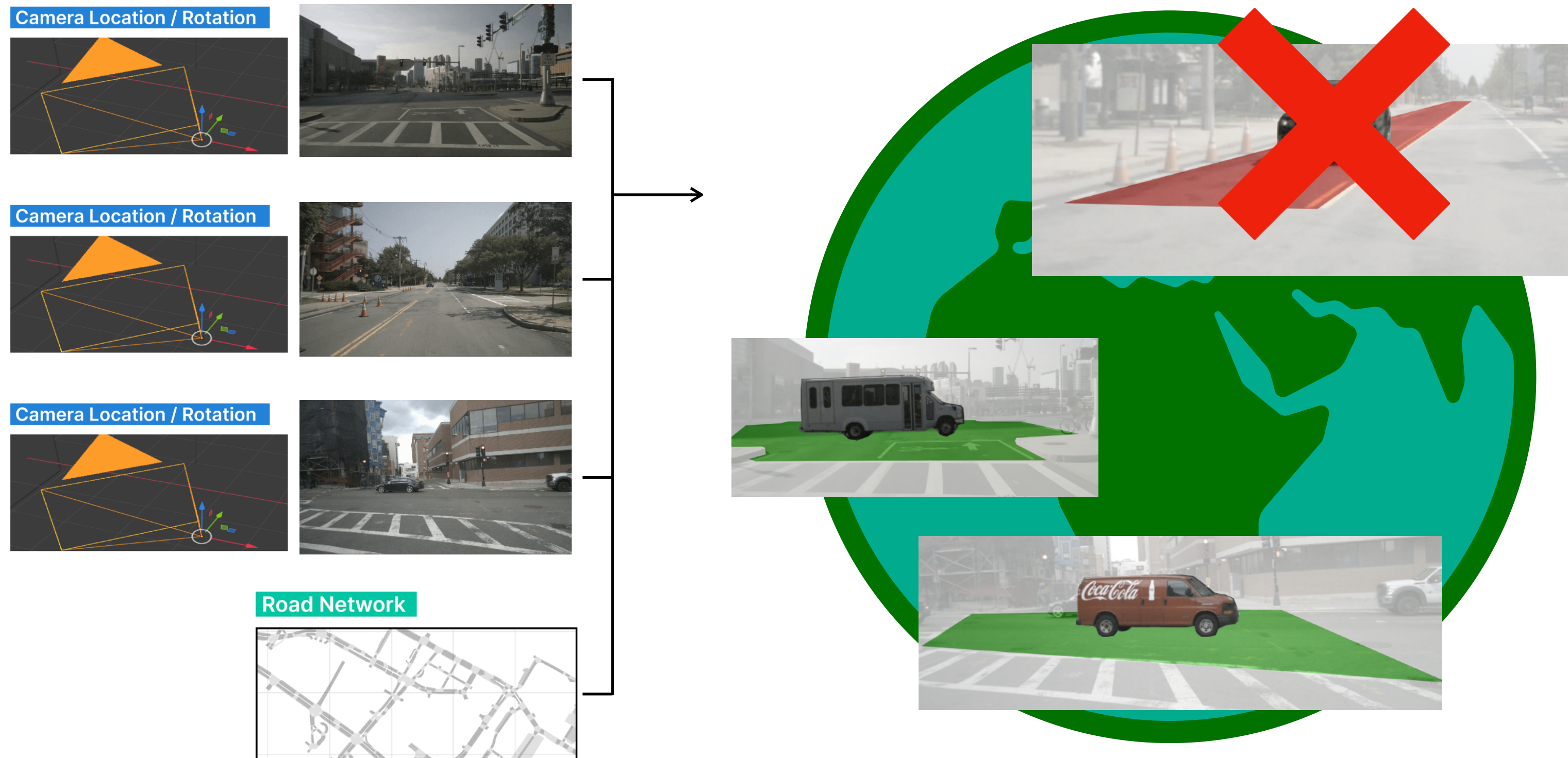
Filter

```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
         contains(i, o))
```

Observe

```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```


Programming Model



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

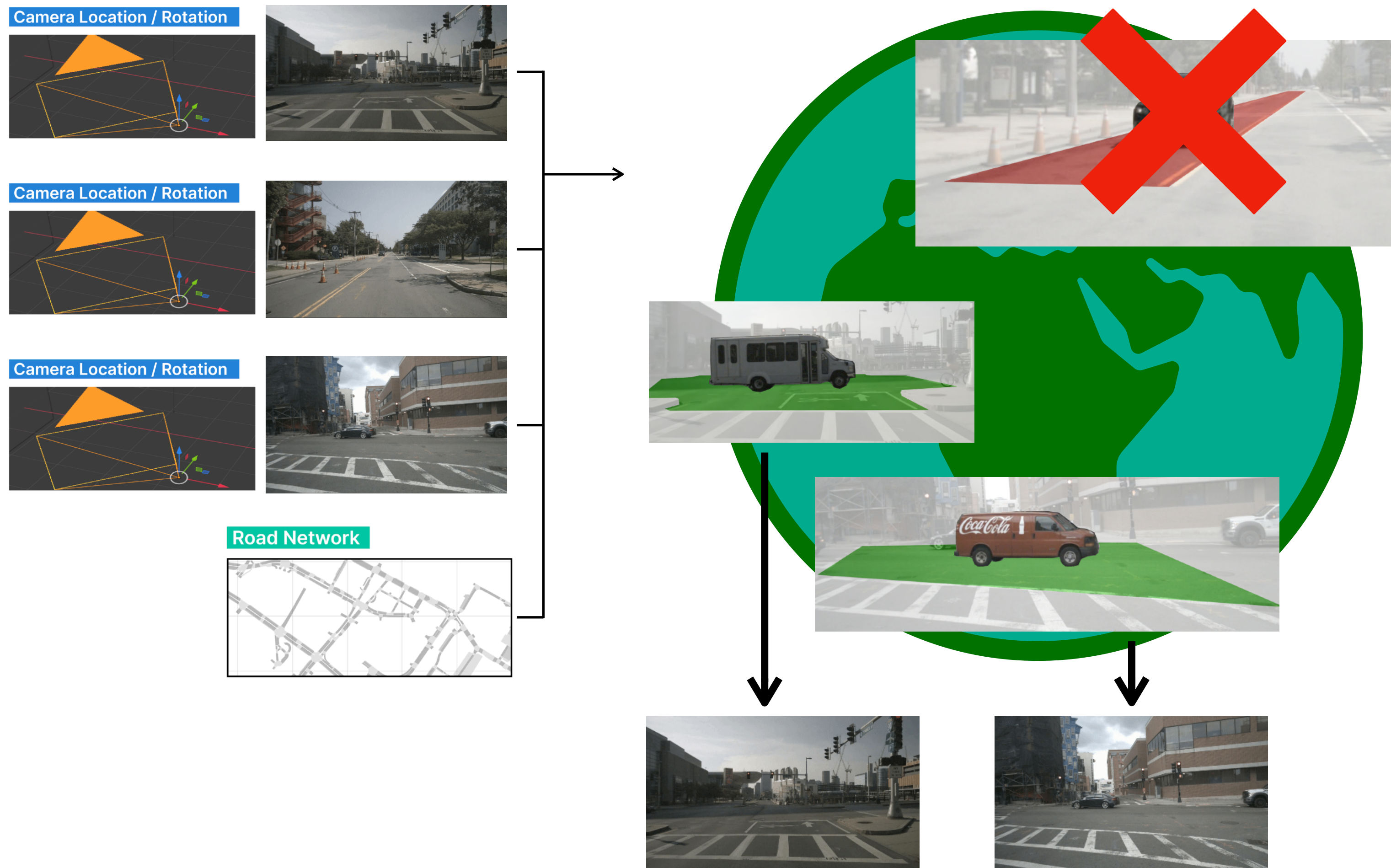
Filter

```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
        contains(i, o))
```

Observe

```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Programming Model



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

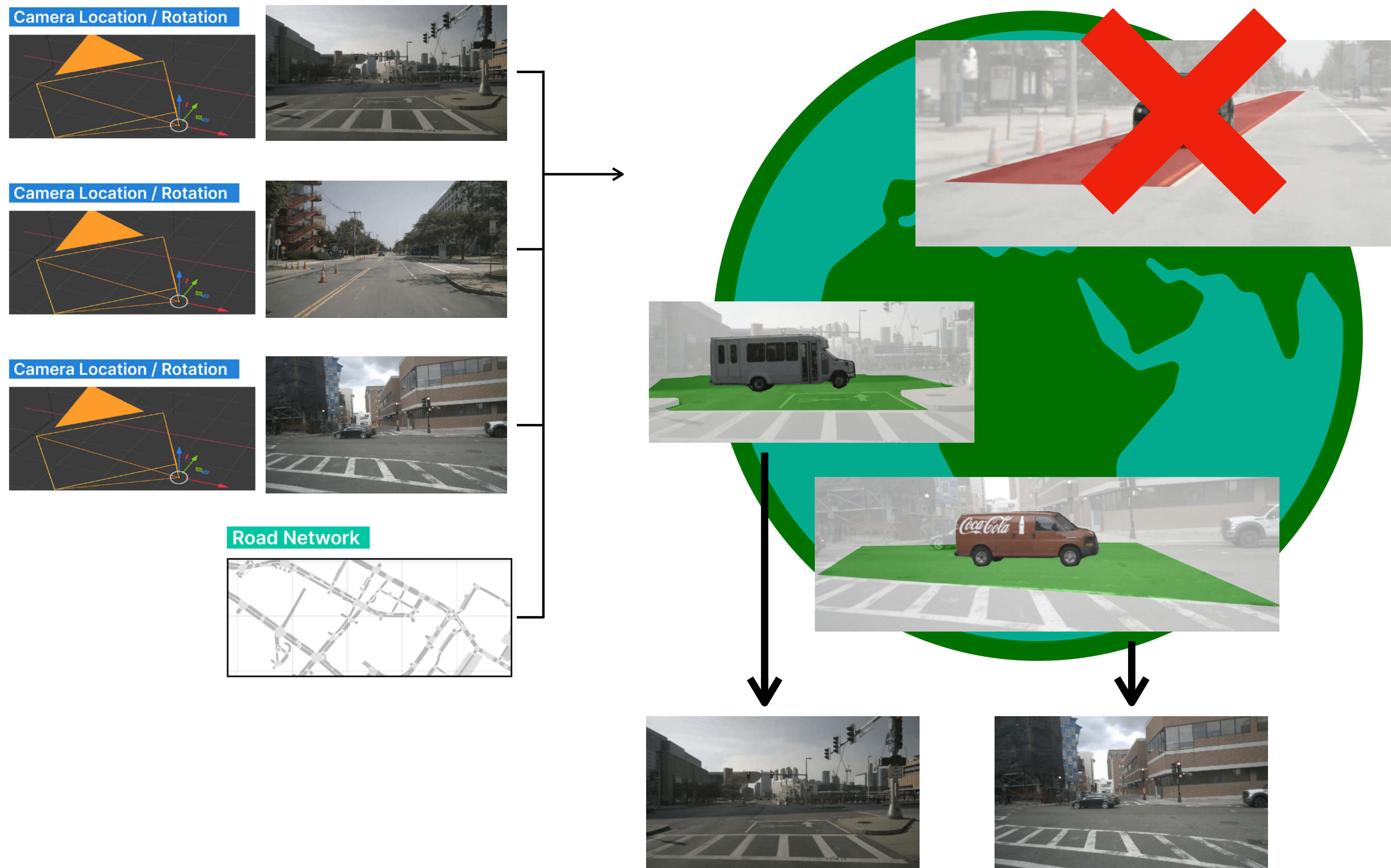
Filter

```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
        contains(i, o))
```

Observe

```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Programming Model



Build

```
w = World()
w.addVideo(video1, camera1)
w.addVideo(video2, camera2)
w.addVideo(video3, camera3)
w.addGeogConstructs(ROADNET_DIR)
```

Filter

```
o = w.object()
i = w.geogConstruct('intersection')
w.filter((o.type=='car') &
        contains(i, o))
```

Observe

```
w.saveVideos('output-dir')
# or
objects = w.getObjects()
```

Video processing execution happens here!

SPATIALYZE

#3 System Optimization

- * Video Processor
- * Optimization Techniques

Video Processor



VIDEO PROCESSOR

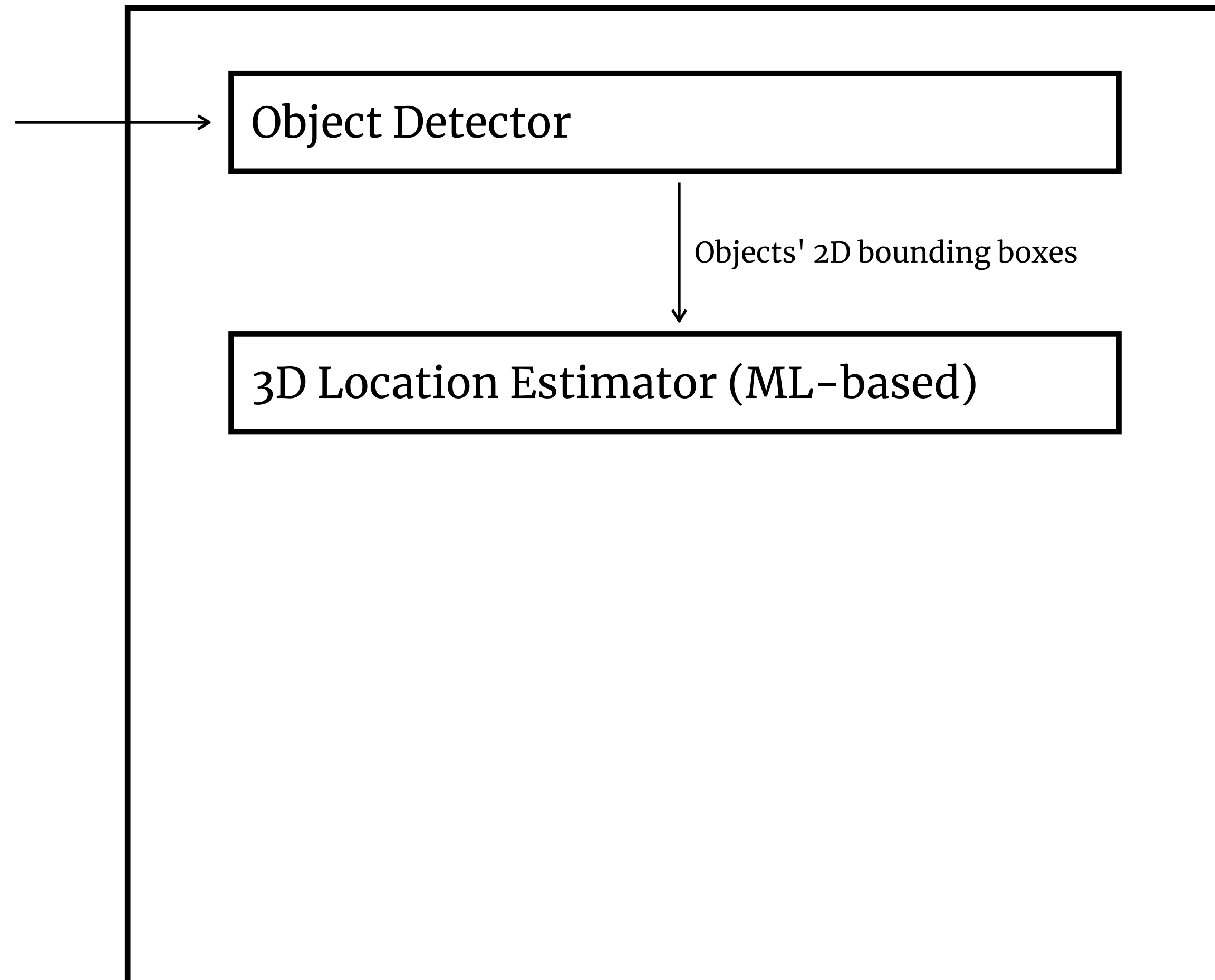


Object Detector

Video Processor



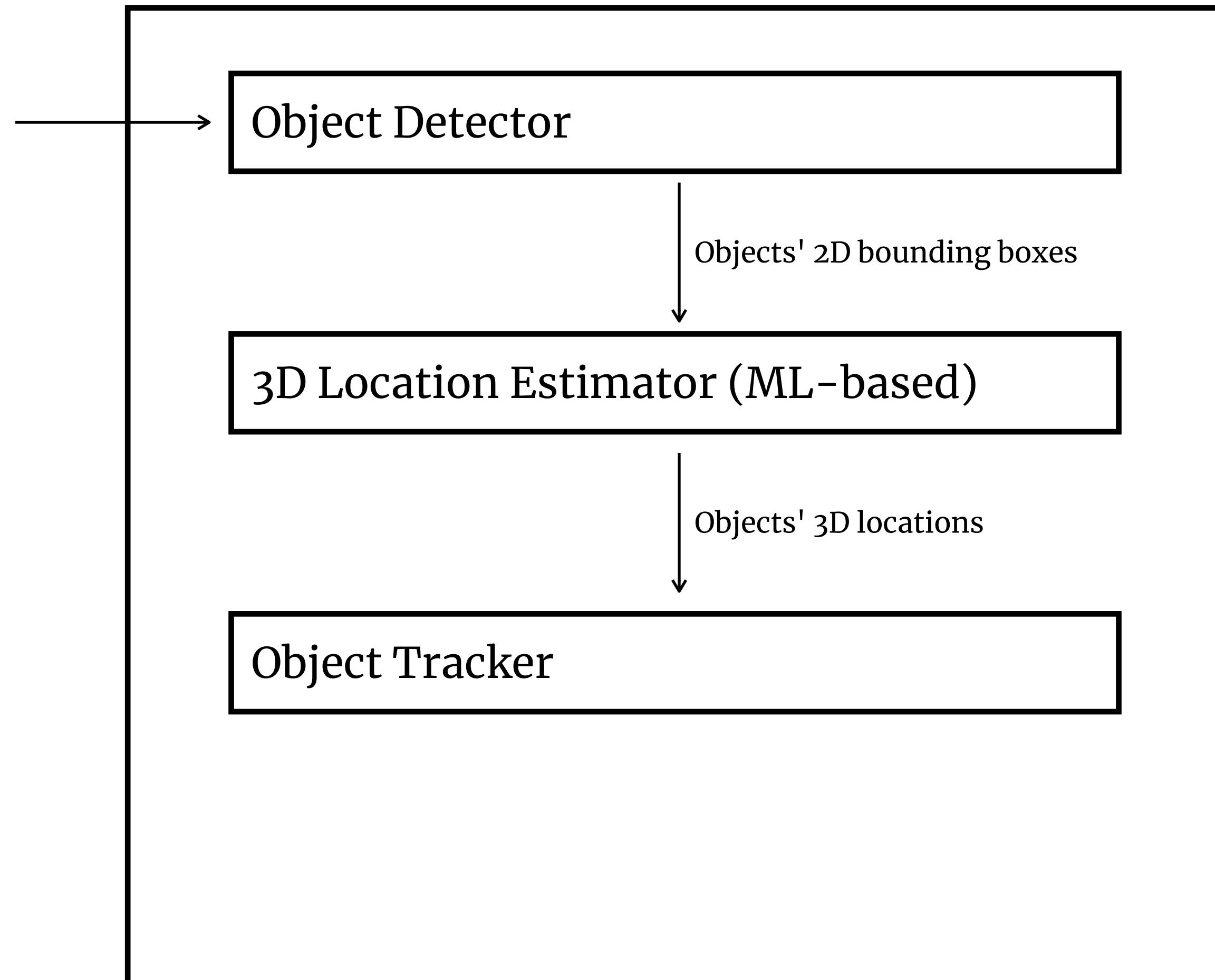
VIDEO PROCESSOR



Video Processor



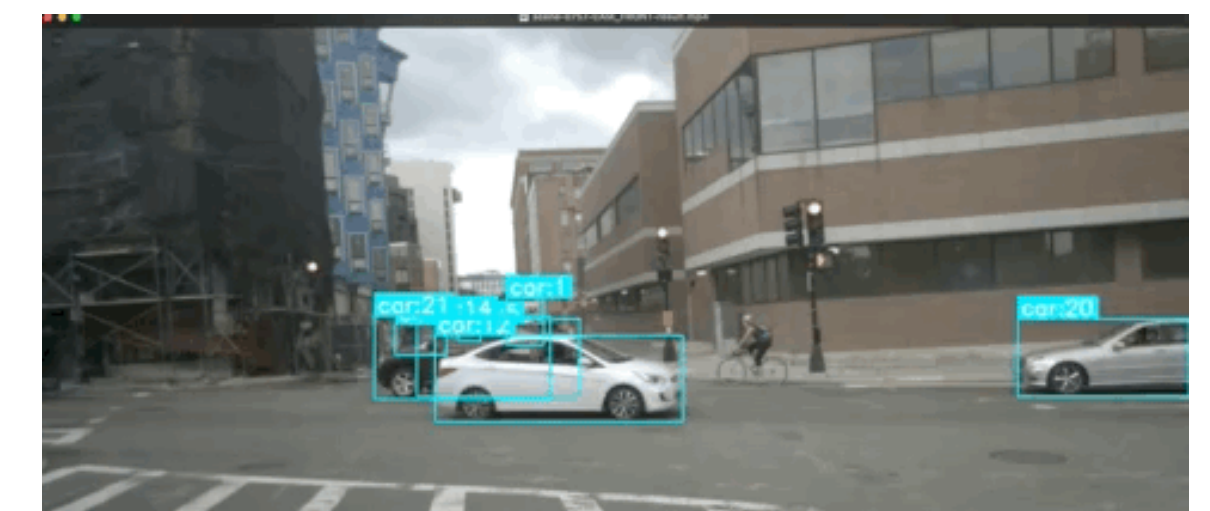
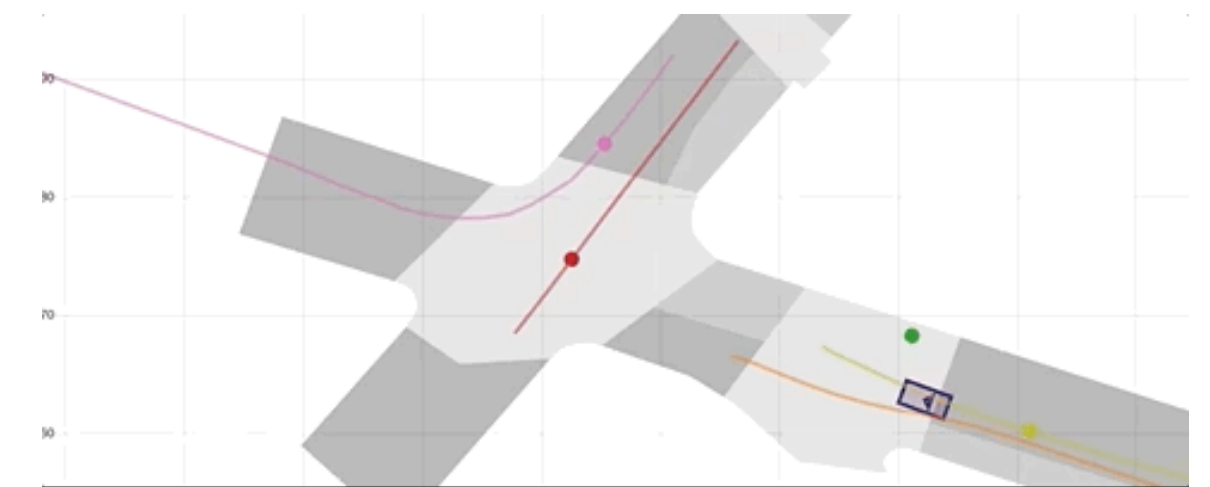
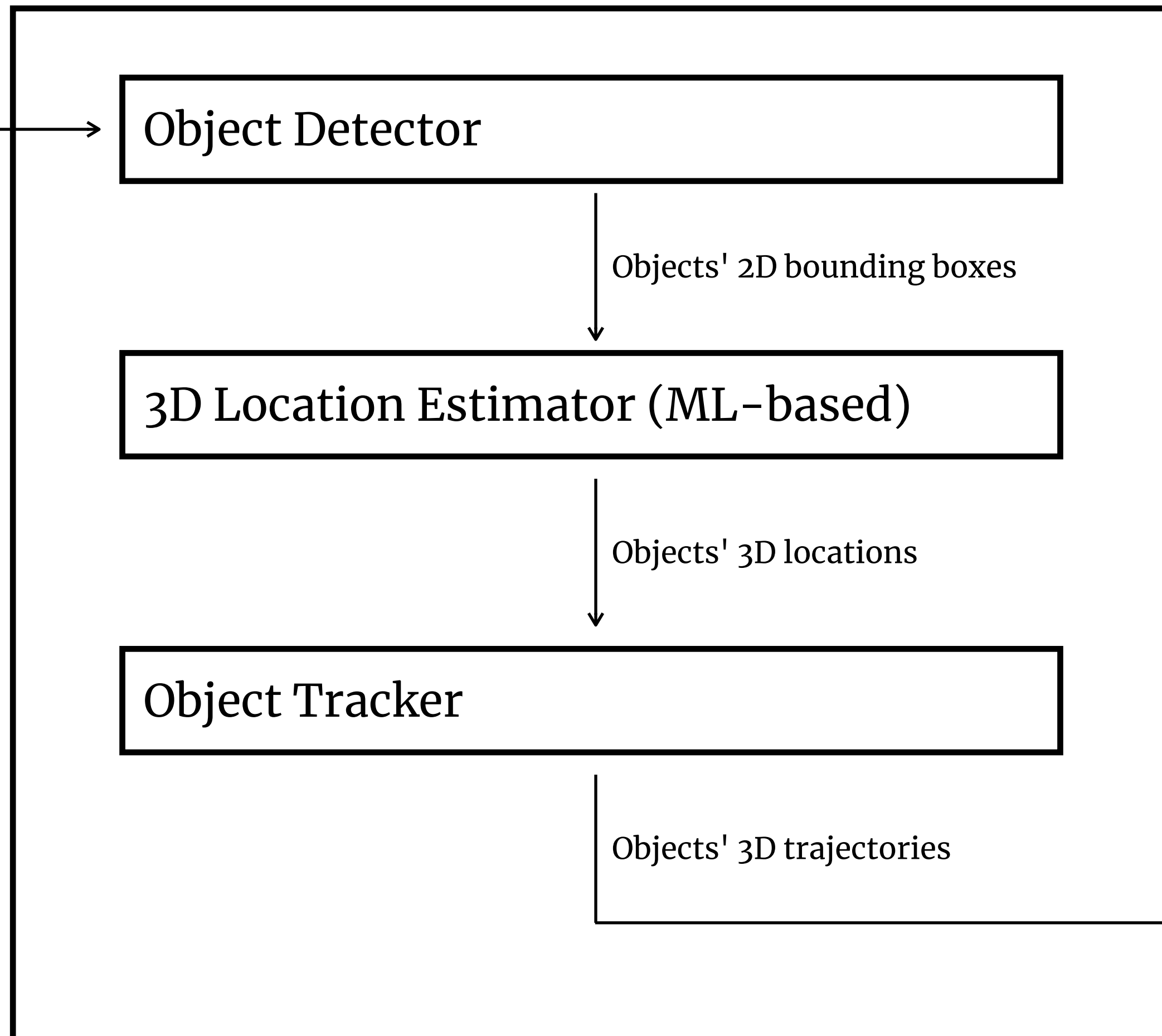
VIDEO PROCESSOR



Video Processor



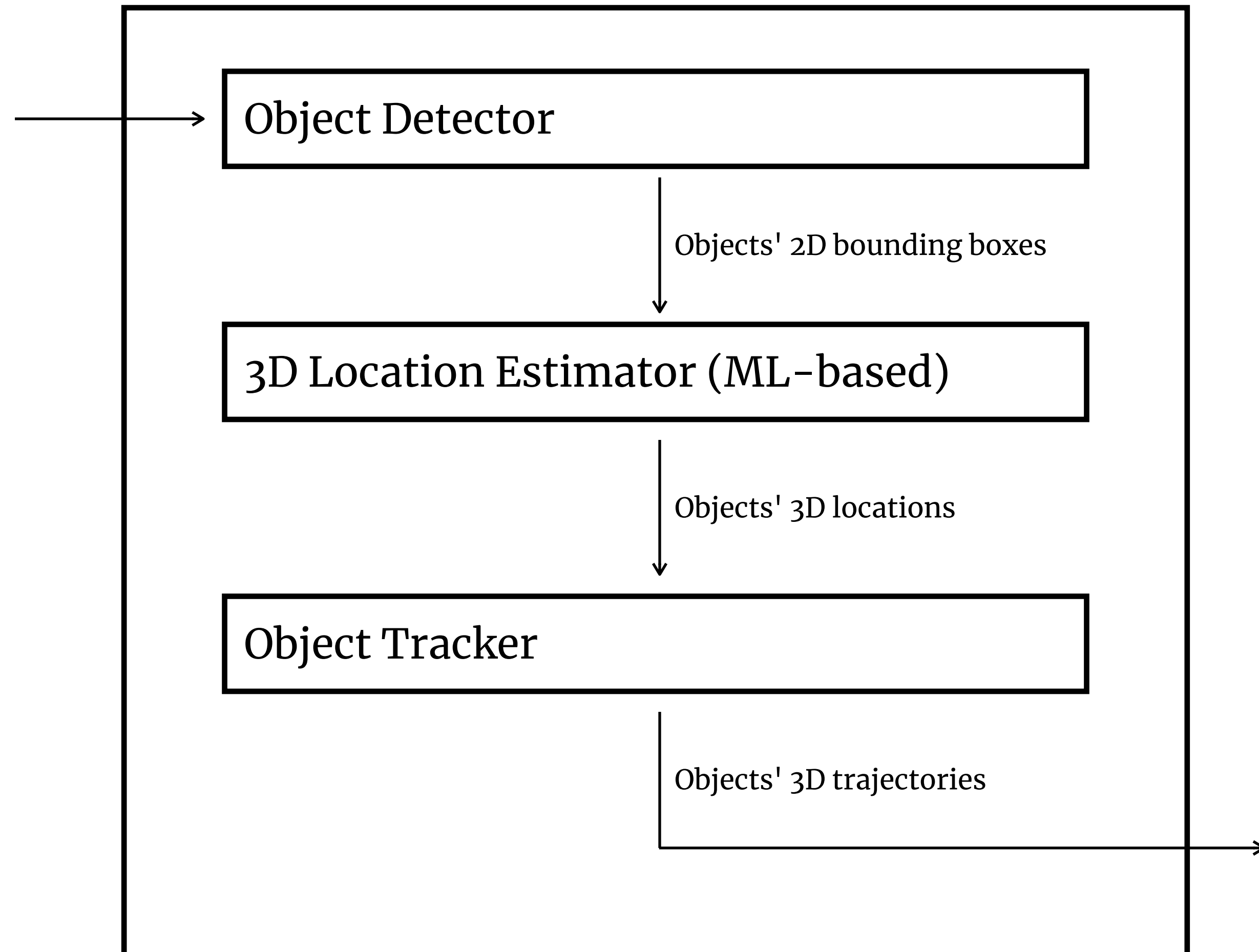
VIDEO PROCESSOR



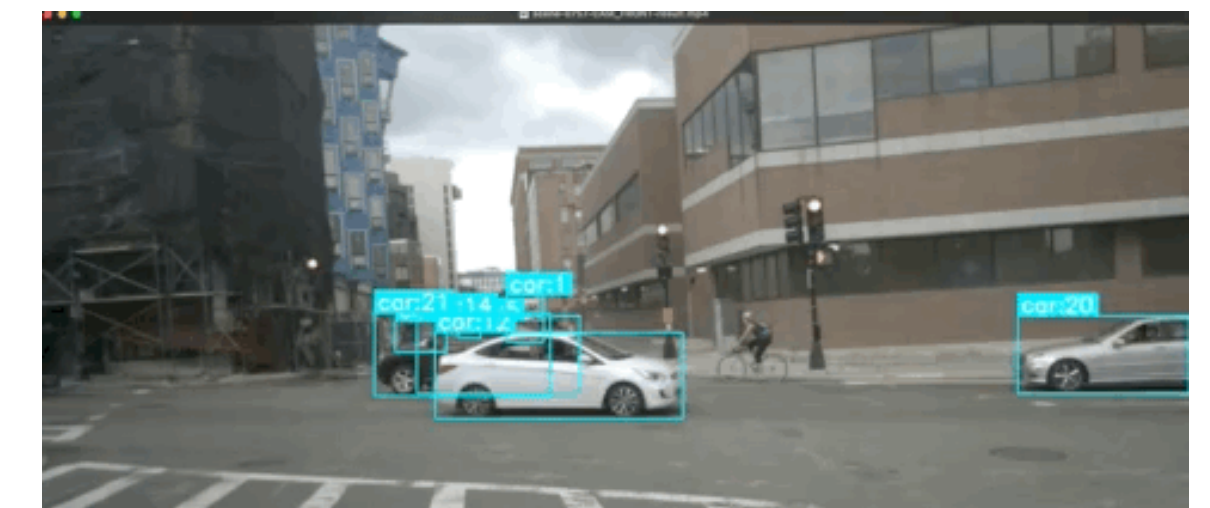
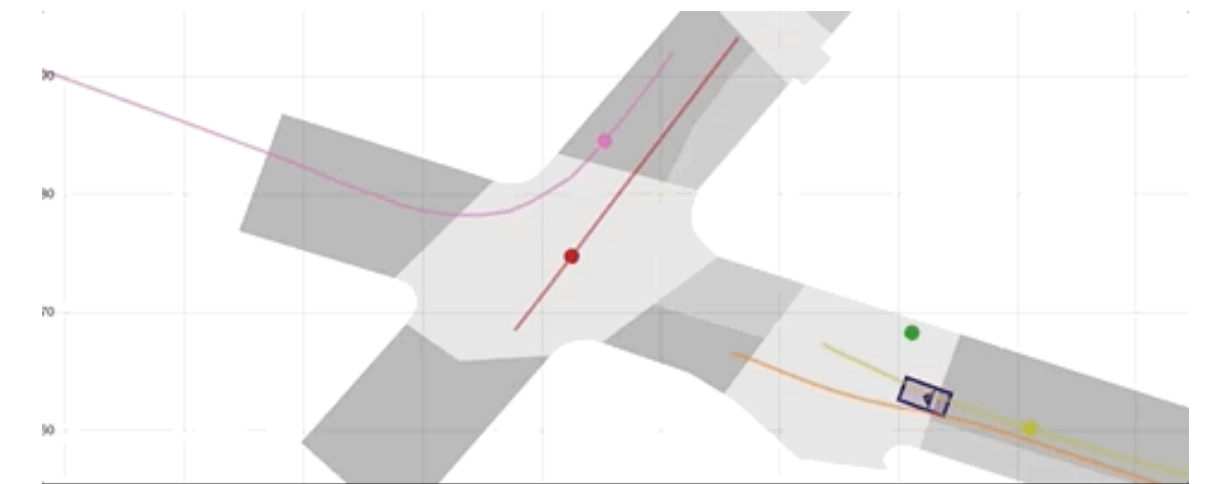
Video Processor



VIDEO PROCESSOR



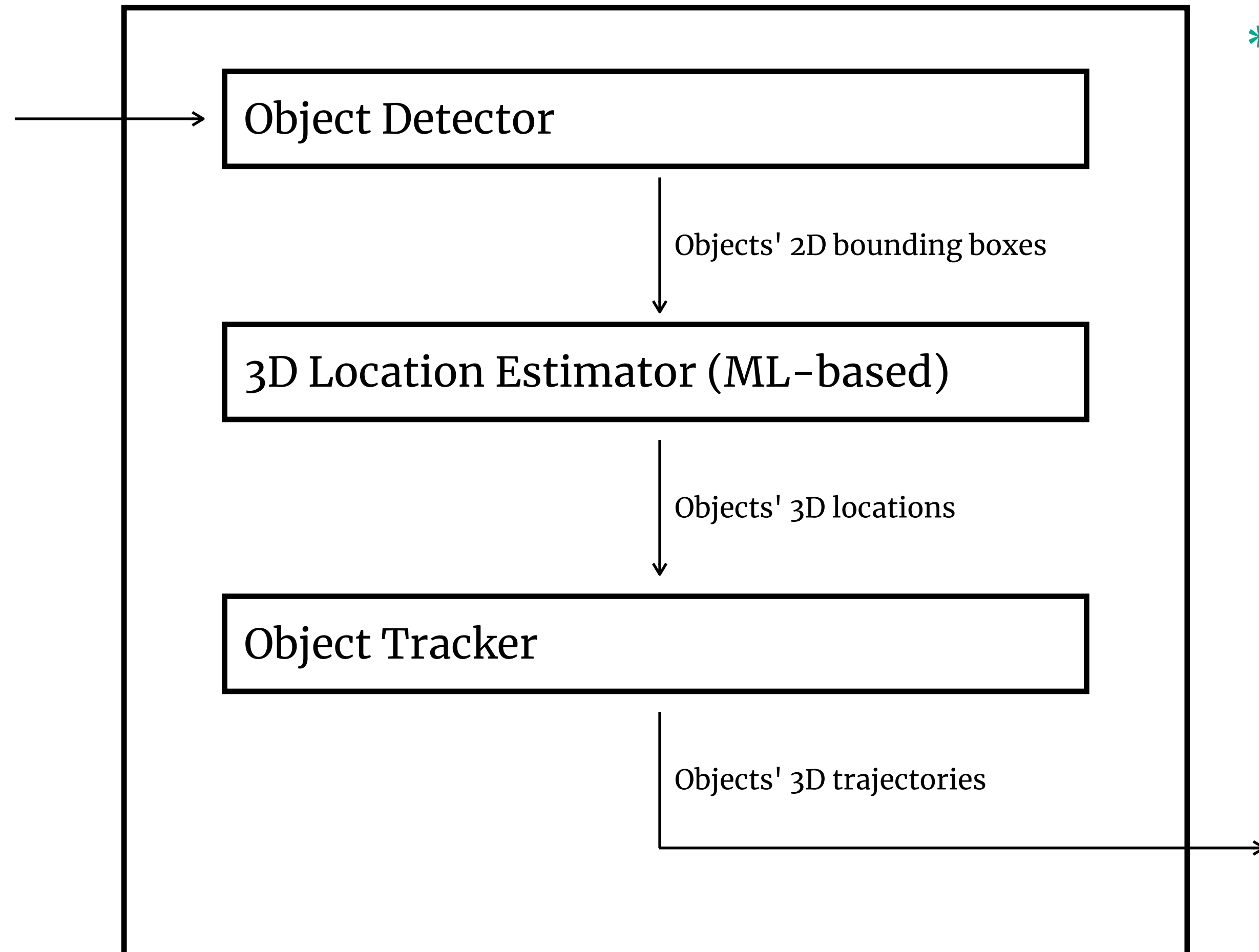
Runtime depends on



Video Processor

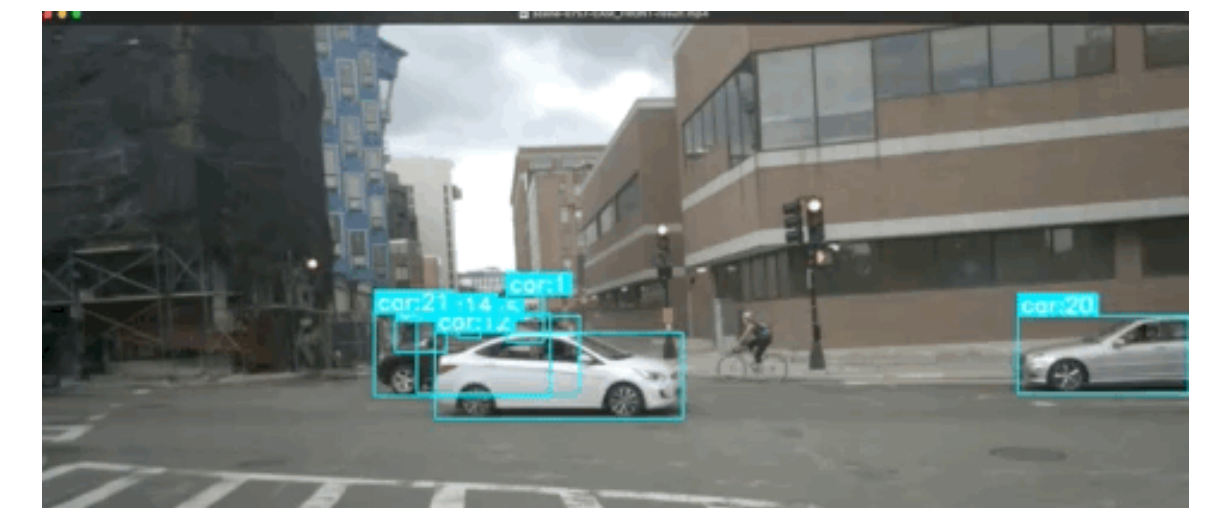
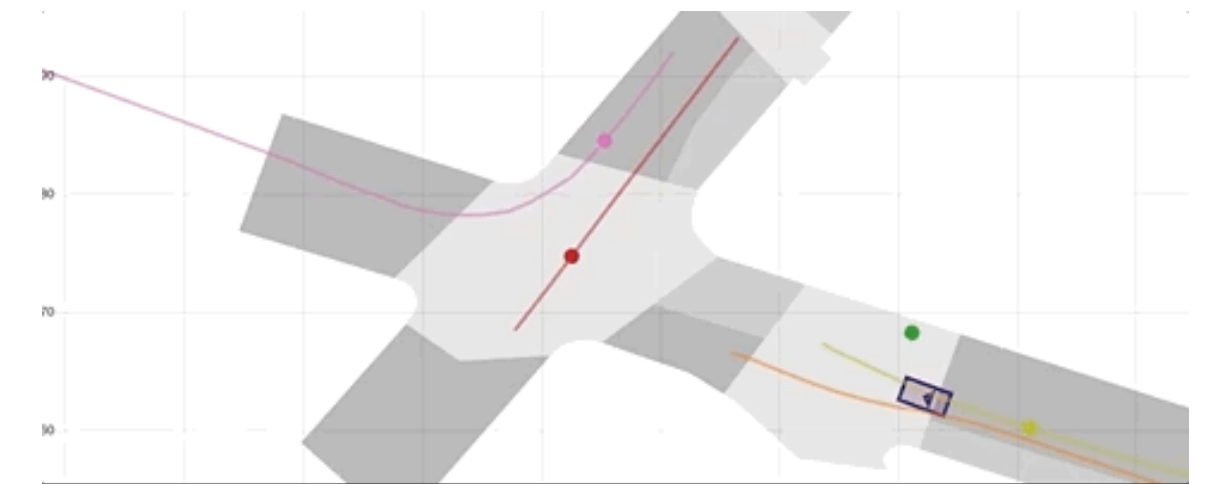


VIDEO PROCESSOR



Runtime depends on

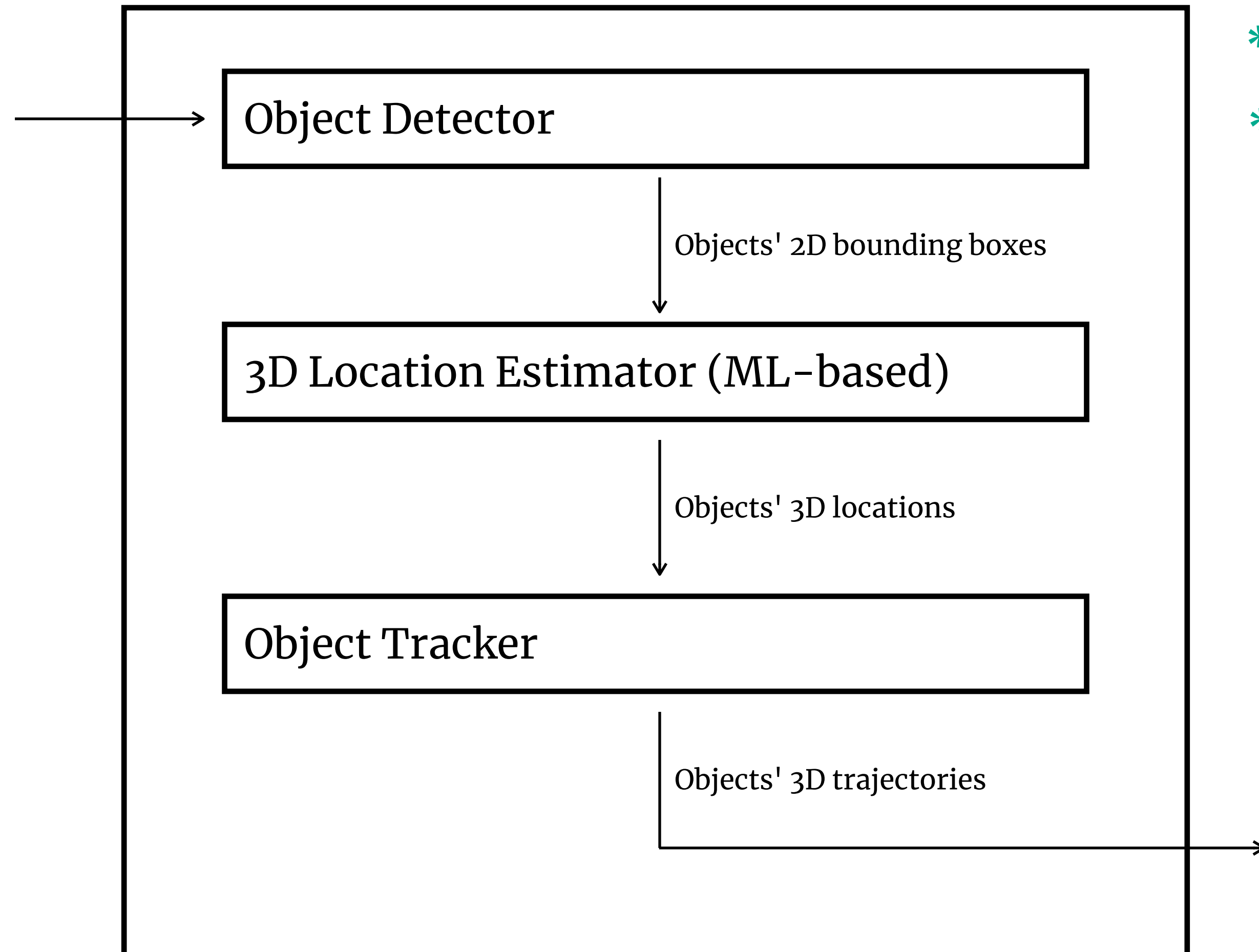
* # of video frames



Video Processor

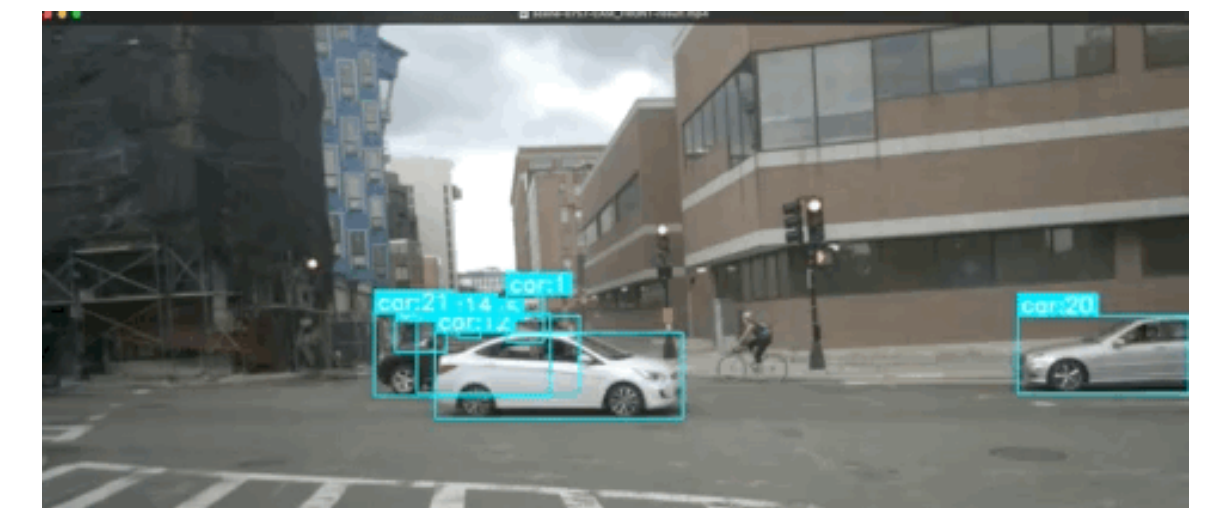
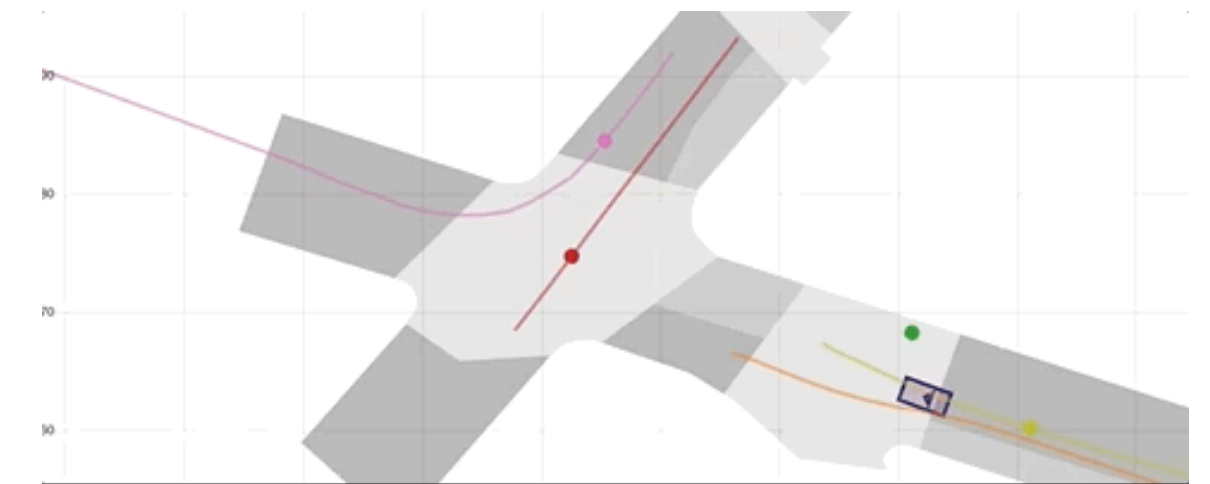


VIDEO PROCESSOR



Runtime depends on

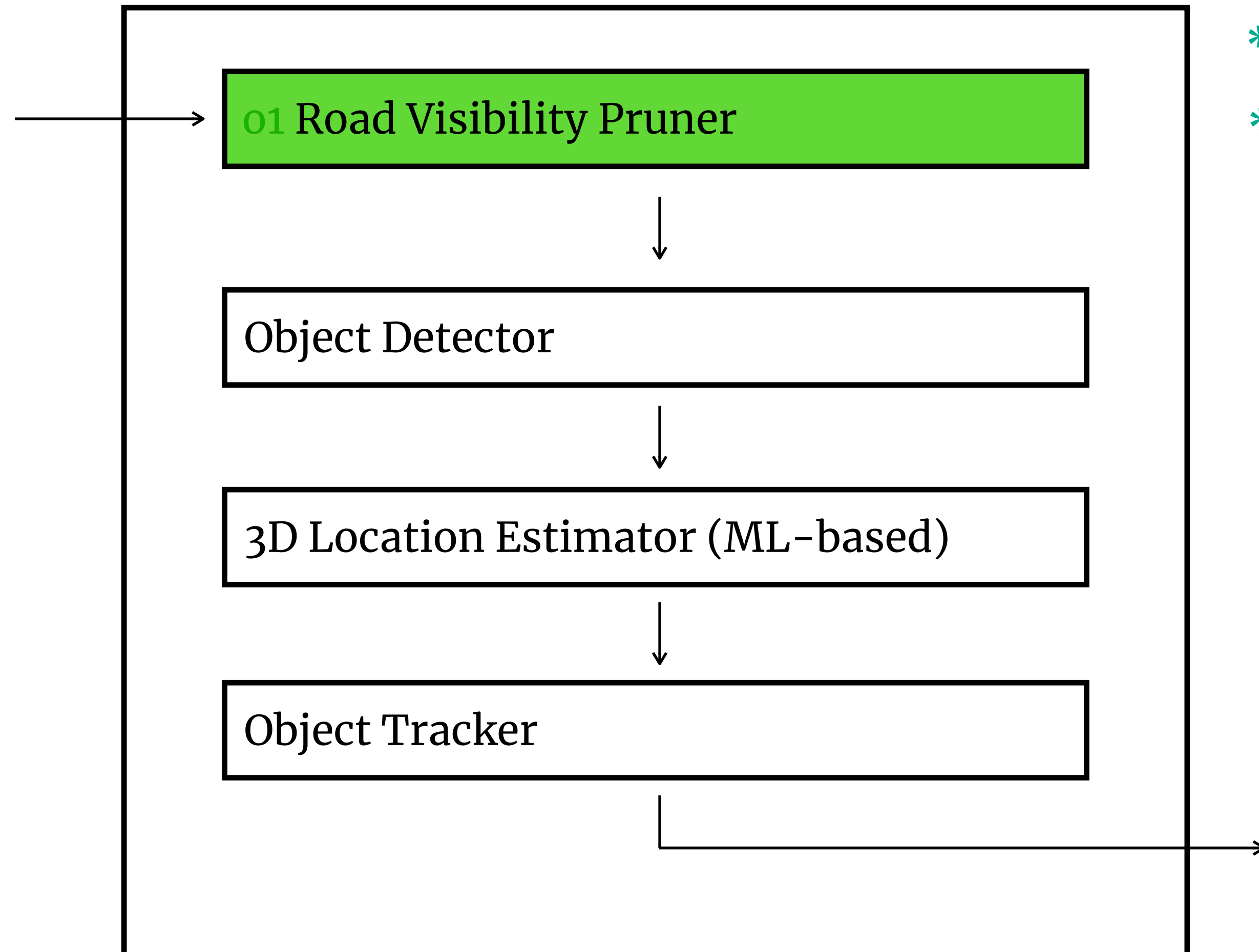
- * # of video frames
- * ML inference time on each video



Video Processor

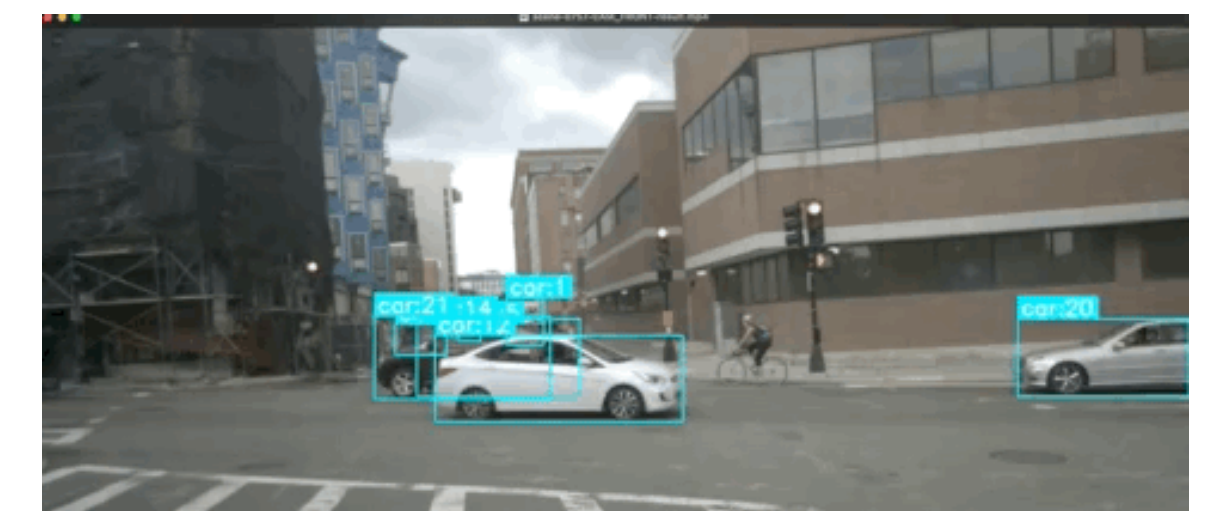
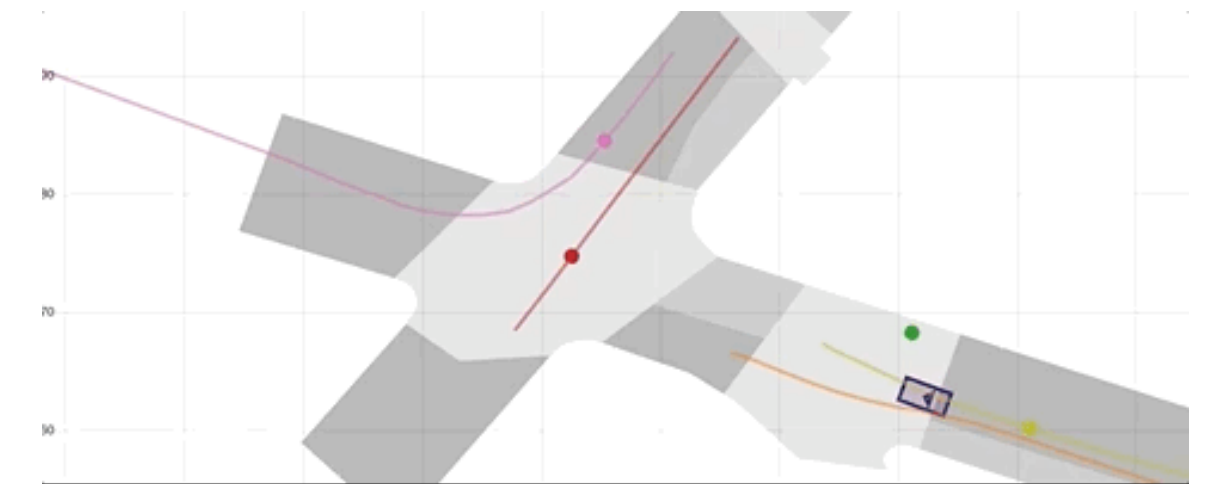


VIDEO PROCESSOR



Runtime depends on

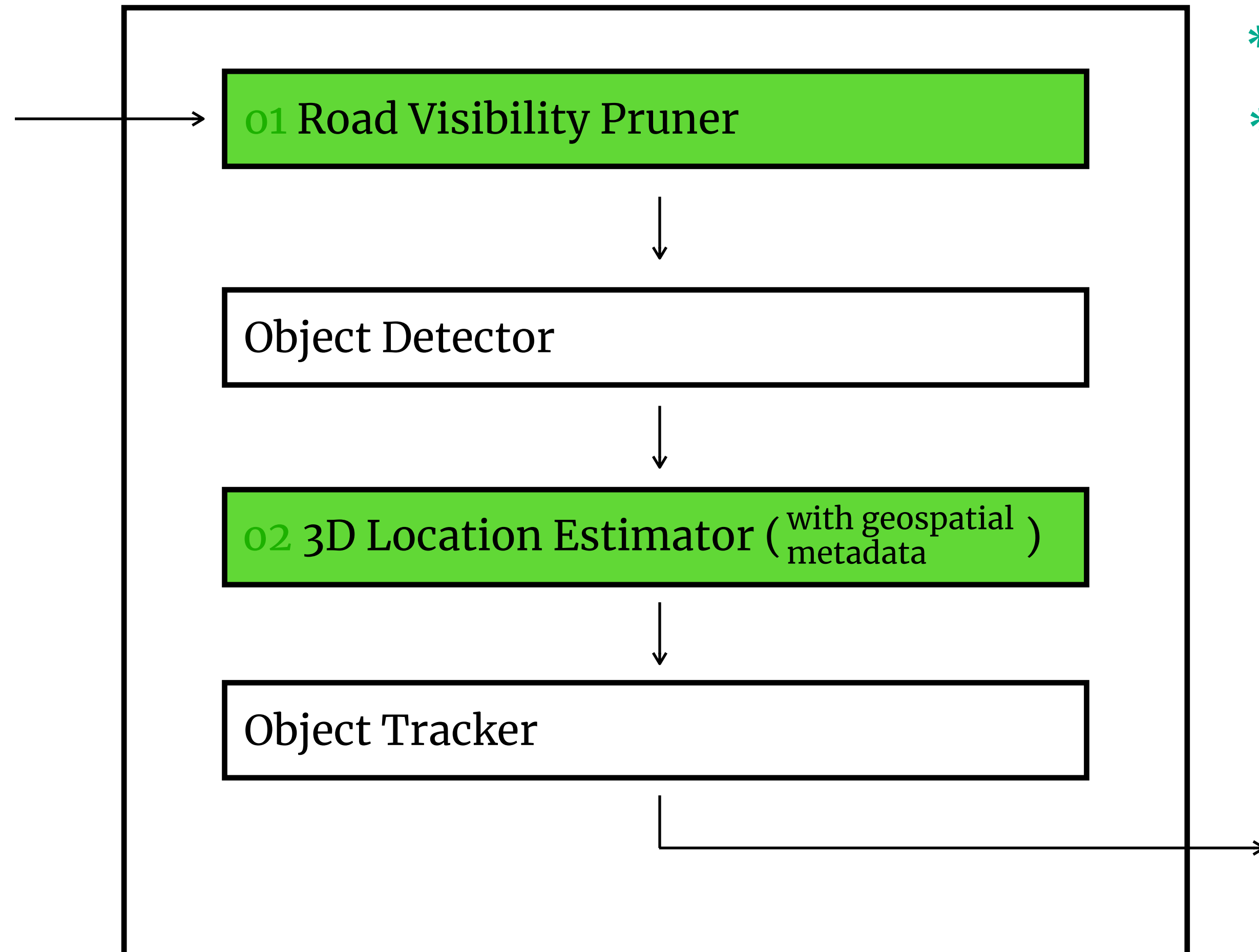
- * # of video frames 01
- * ML inference time on each video



Video Processor

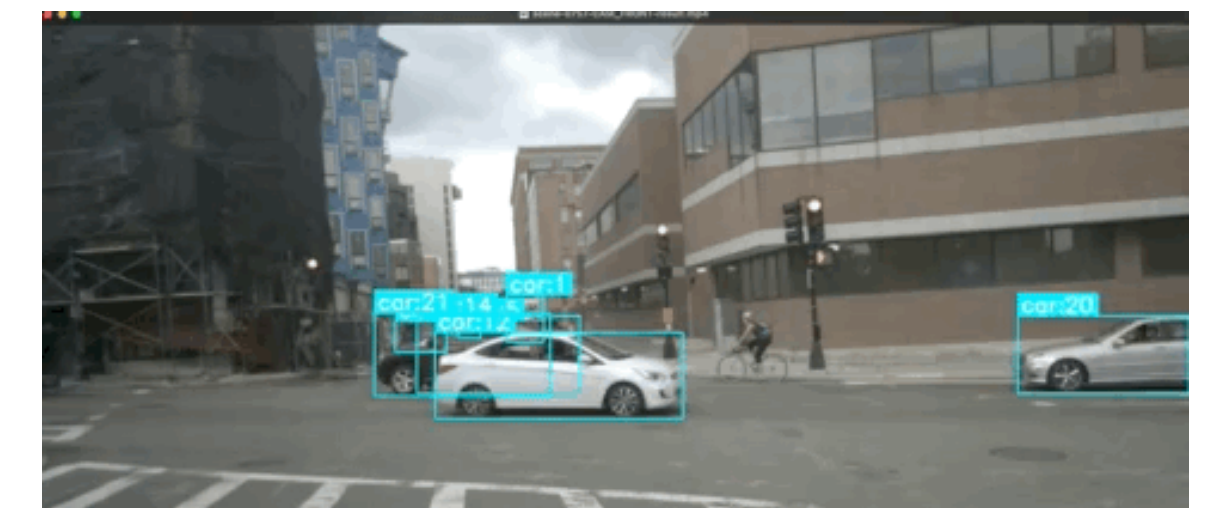
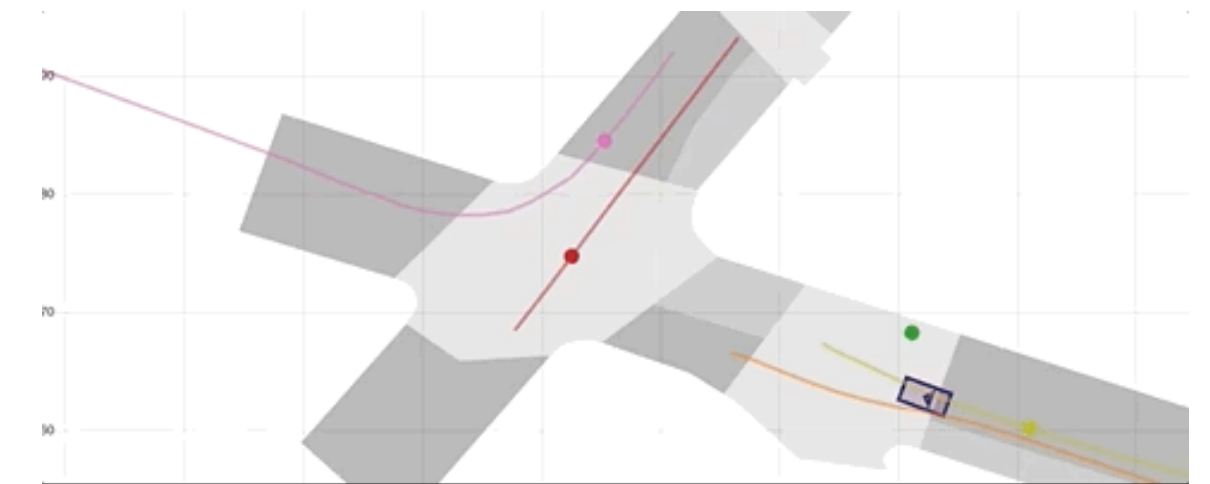


VIDEO PROCESSOR



Runtime depends on

- * # of video frames 01
- * ML inference time on each video 02



Optimization Techniques 01 Road Visibility Pruner

QUERY

```
contains(intersection, car)
```



IMPLICATION

The car must be at an intersection



The intersection must be in the camera's view

Optimization Techniques 01 Road Visibility Pruner

QUERY

```
contains(intersection, car)
```



IMPLICATION

The car must be at an intersection



The intersection must be in the camera's view

Optimization Techniques

01 Road Visibility Pruner

QUERY

```
contains(intersection, car)
```



IMPLICATION

The car must be at an intersection



The intersection must be in the camera's view

Optimization Techniques

01 Road Visibility Pruner

QUERY

```
contains(intersection, car)
```



IMPLICATION

The **car** must be at an **intersection**



The intersection must be in the camera's view

Optimization Techniques

01 Road Visibility Pruner

QUERY

```
contains(intersection, car)
```



IMPLICATION

The **car** must be at an **intersection**



The **intersection** must be in the **camera's view**

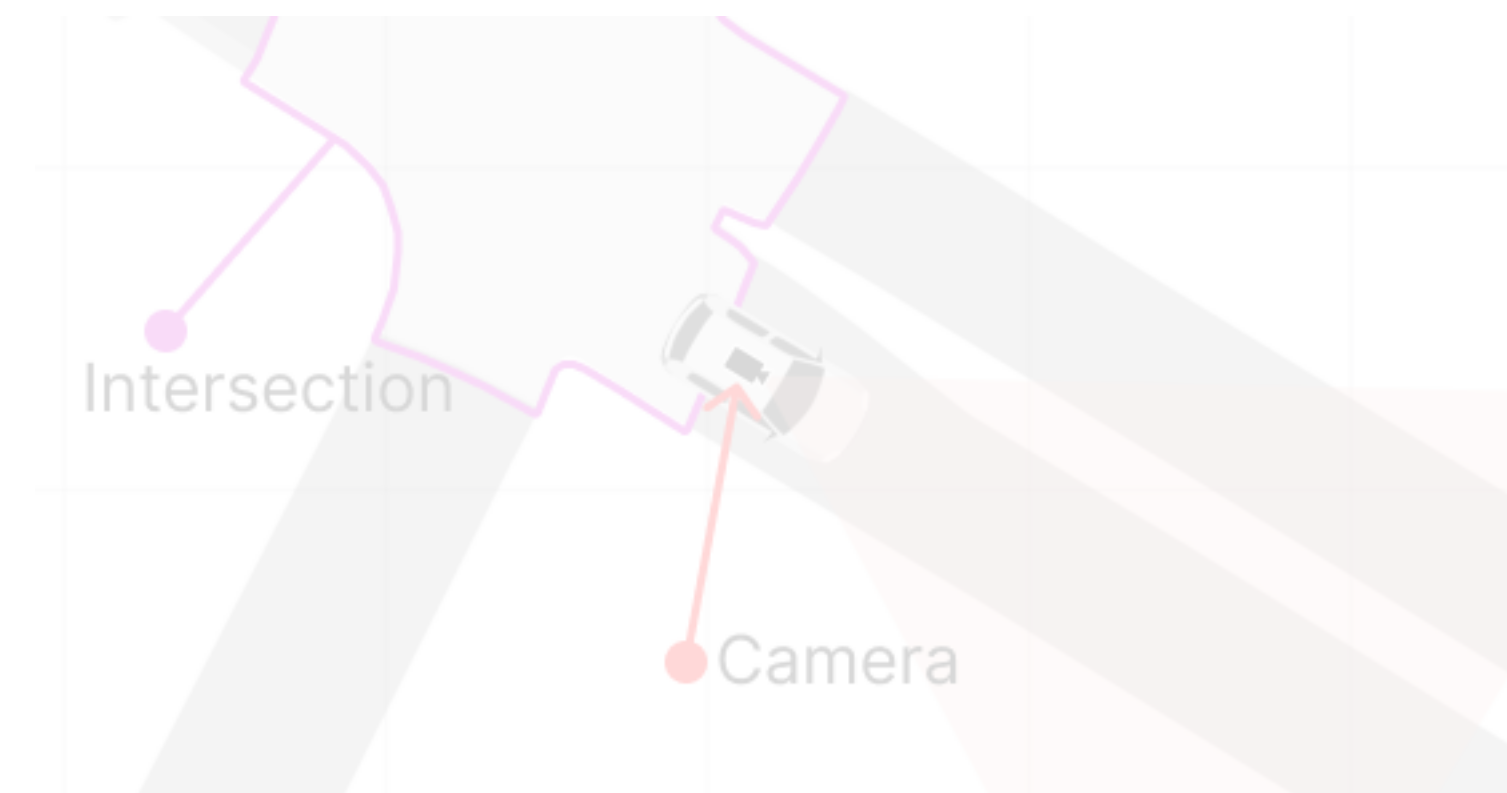
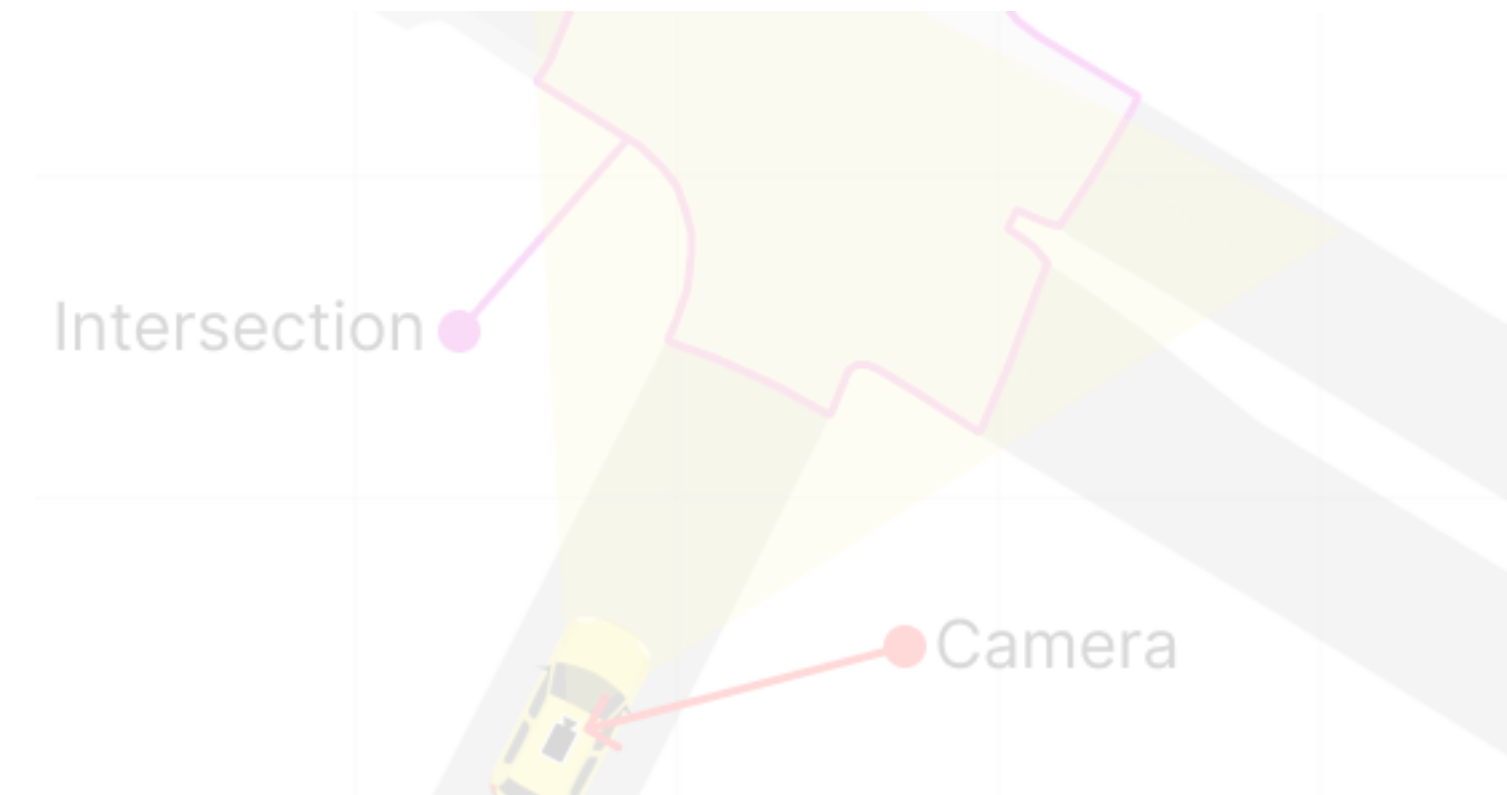
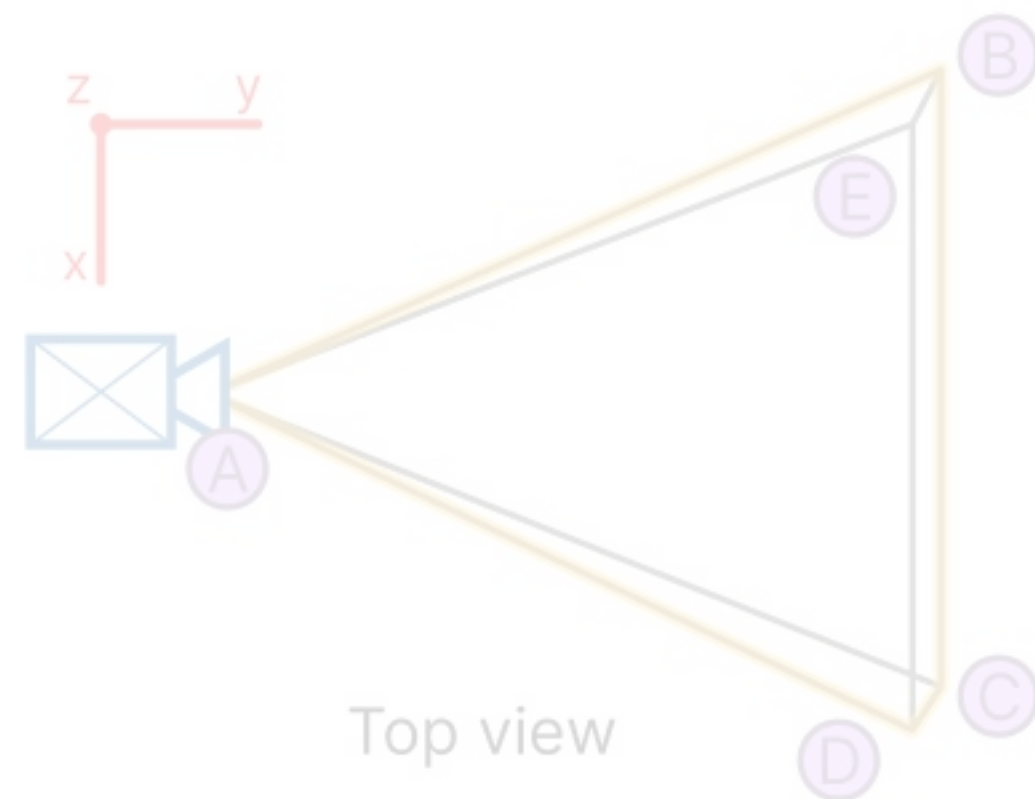
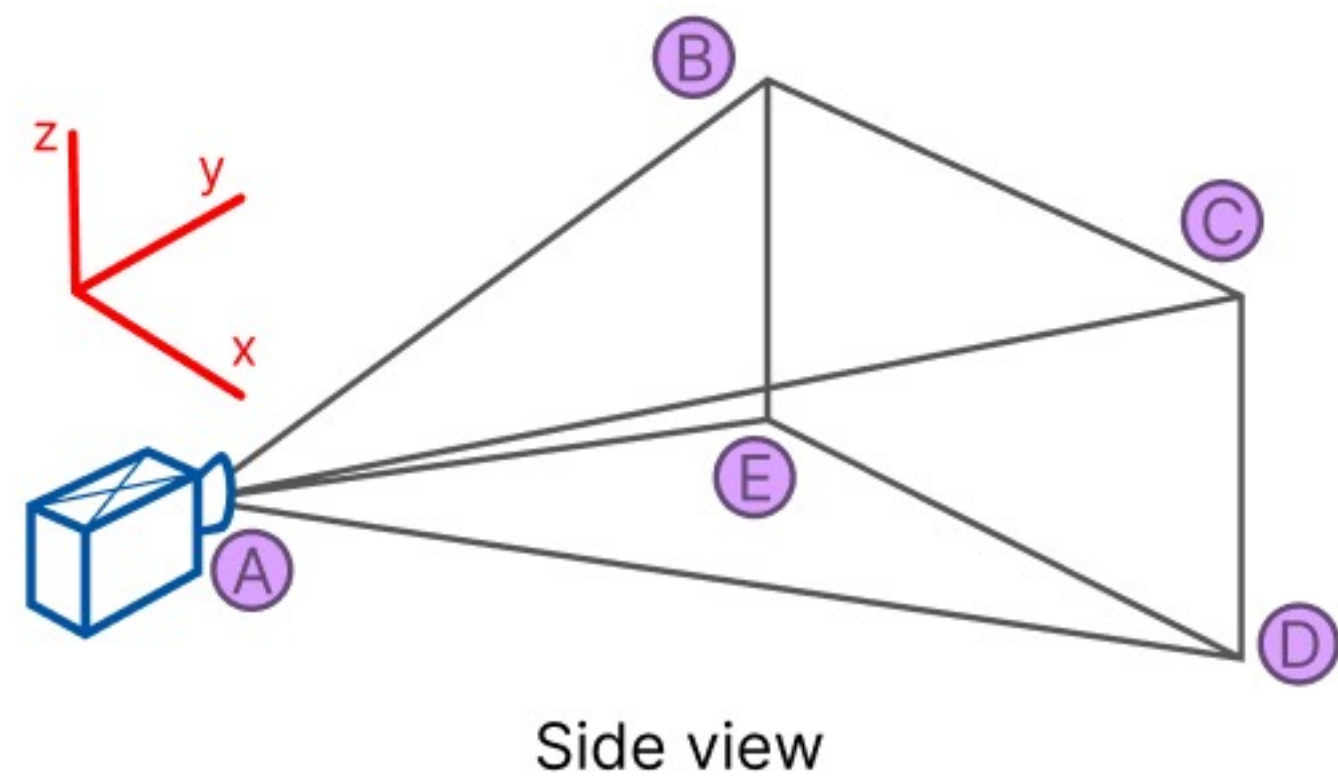
Optimization Techniques

01 Road Visibility Pruner

"**Geographic constructs'** visibility is a proxy for **objects'** visibility"

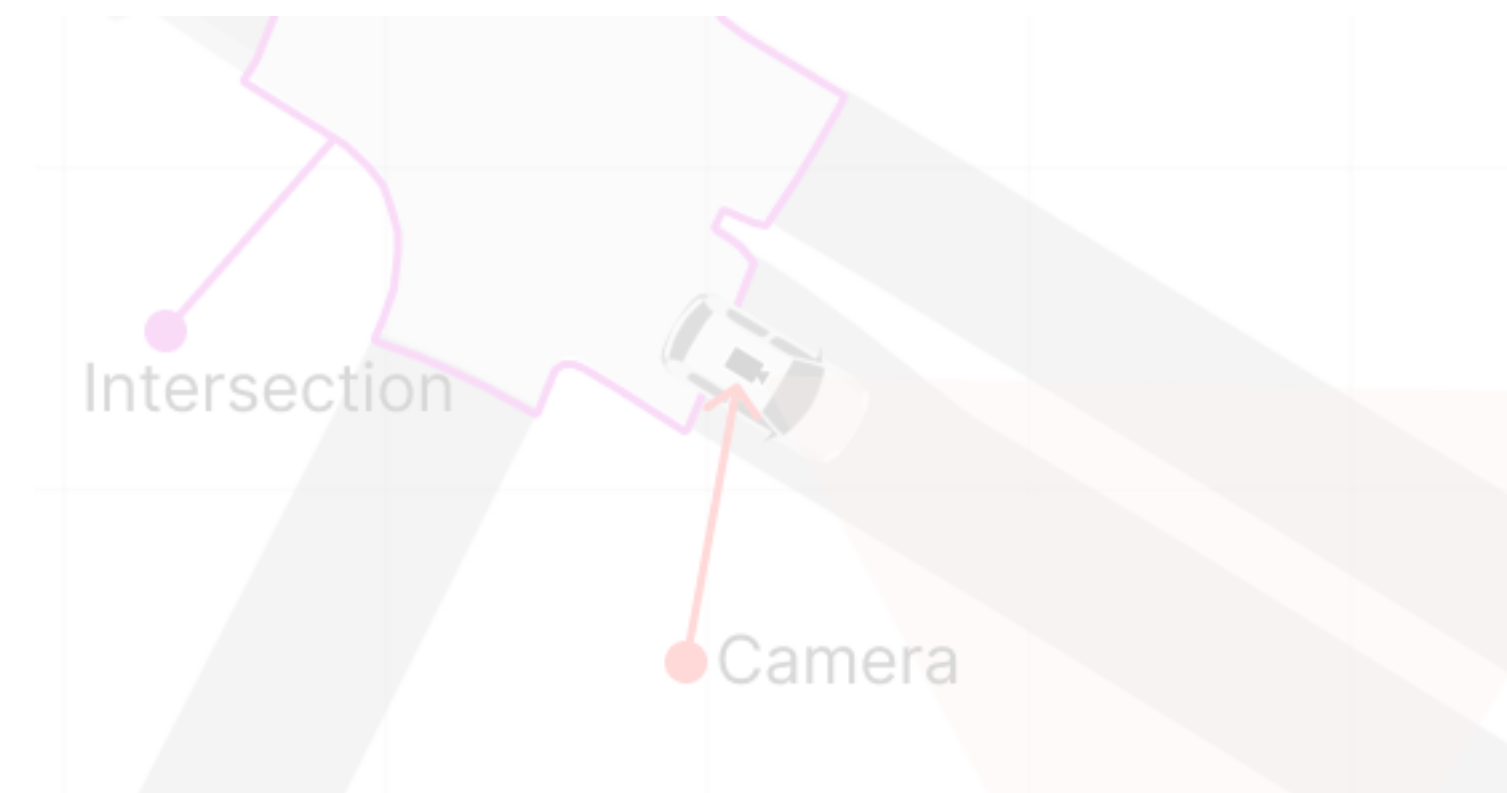
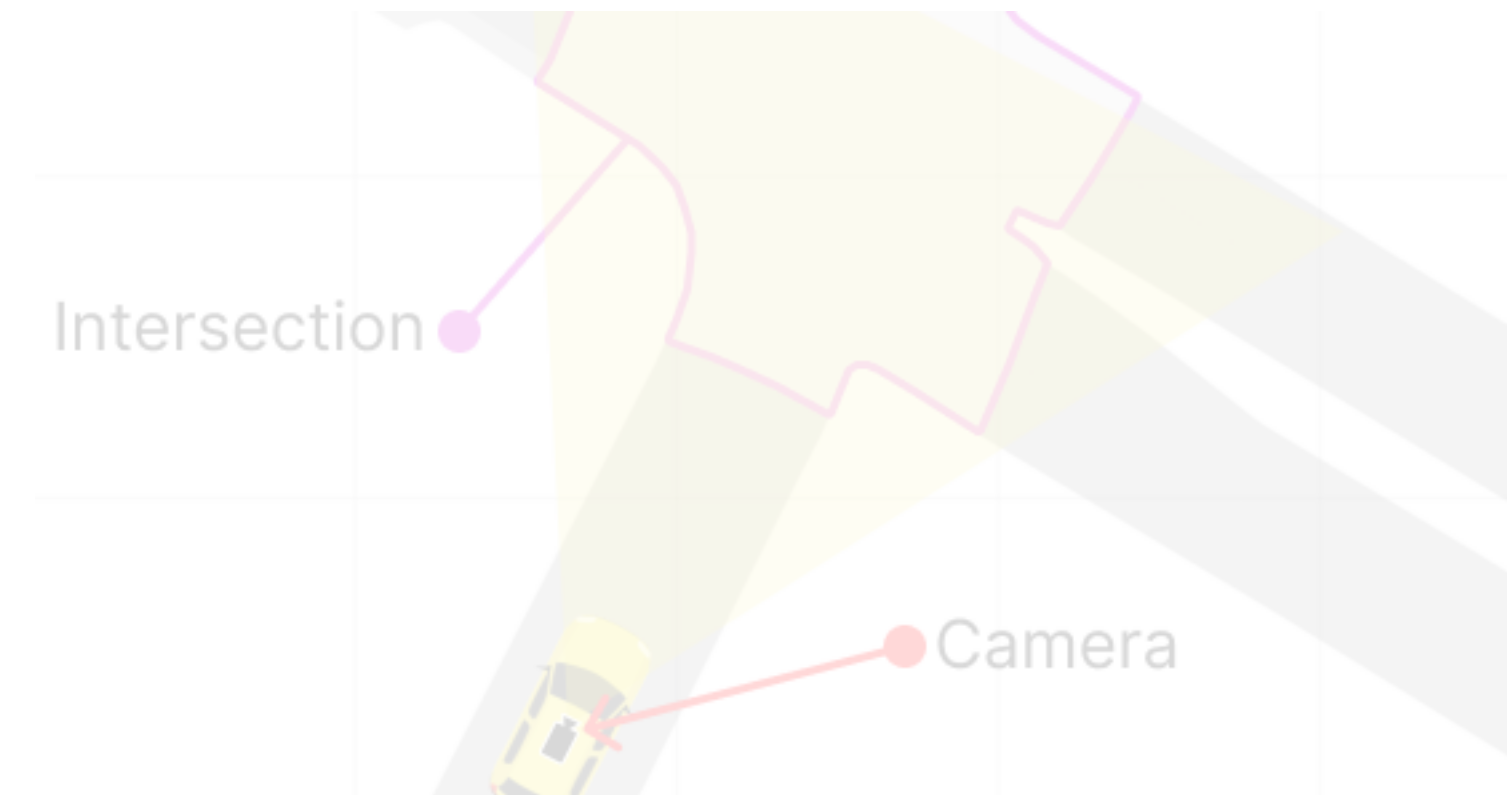
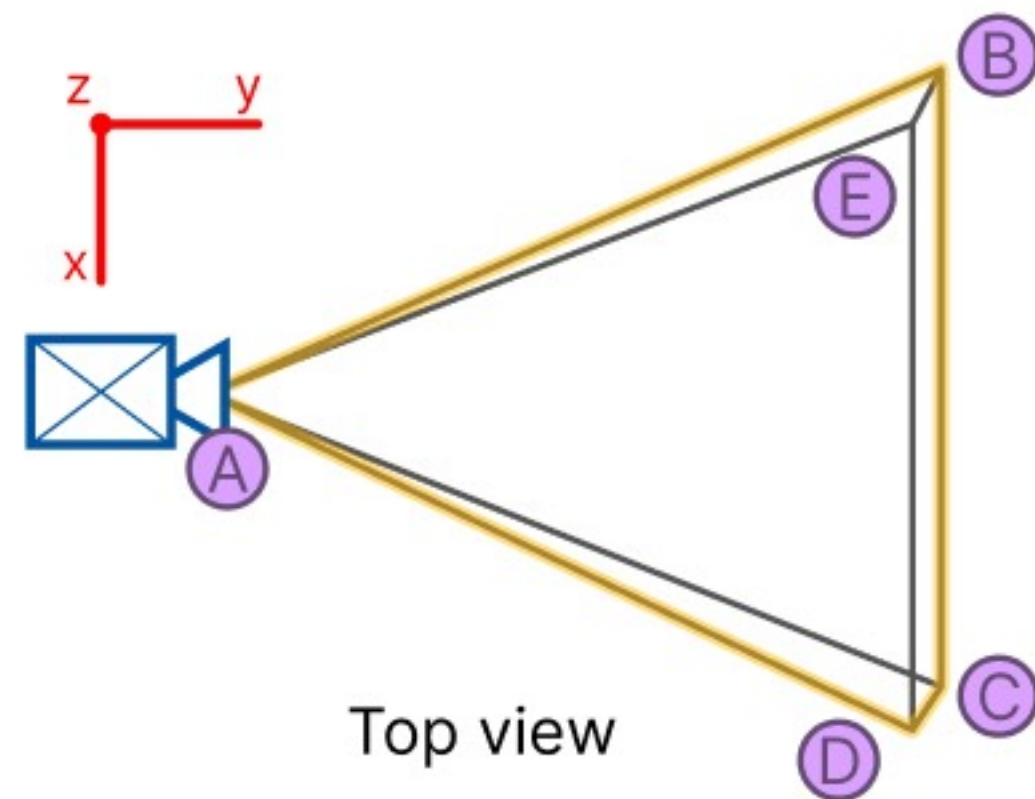
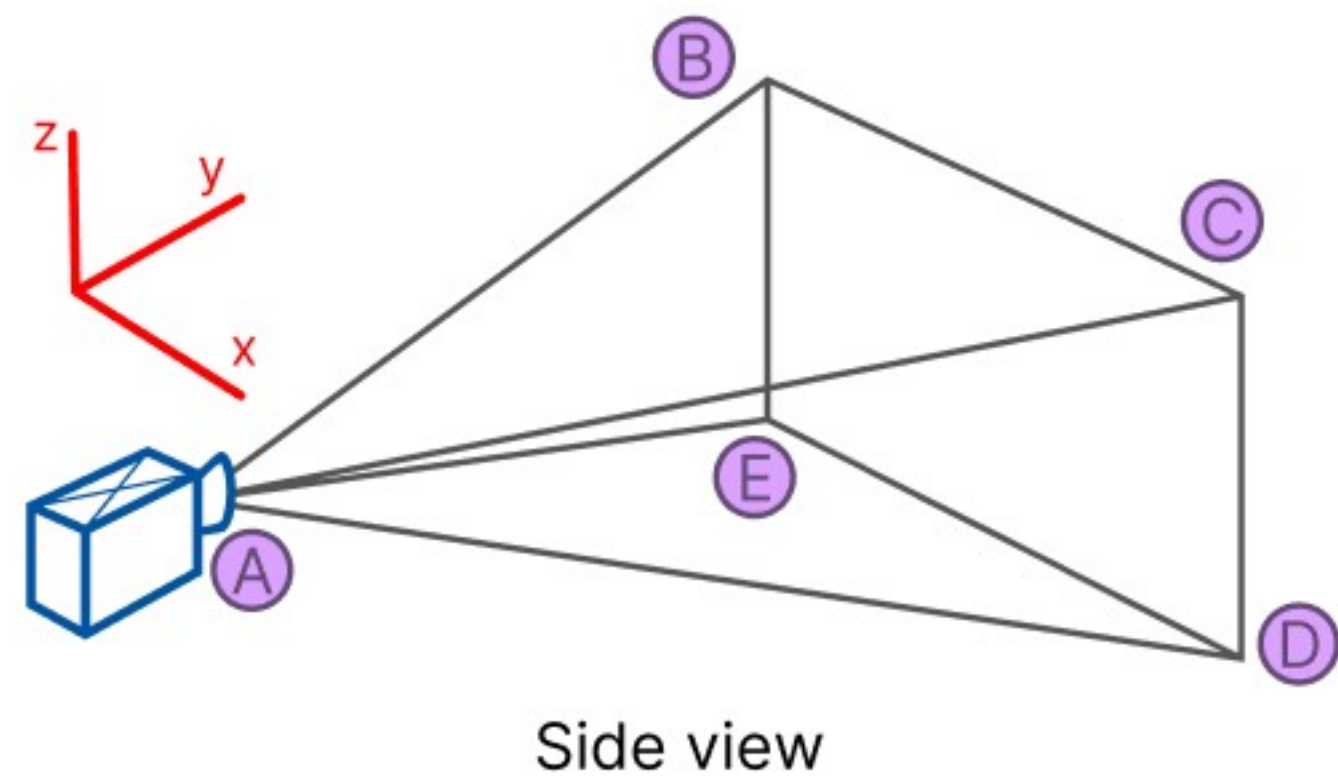
Optimization Techniques 01 Road Visibility Pruner

"Geographic constructs' visibility is a proxy for objects' visibility"



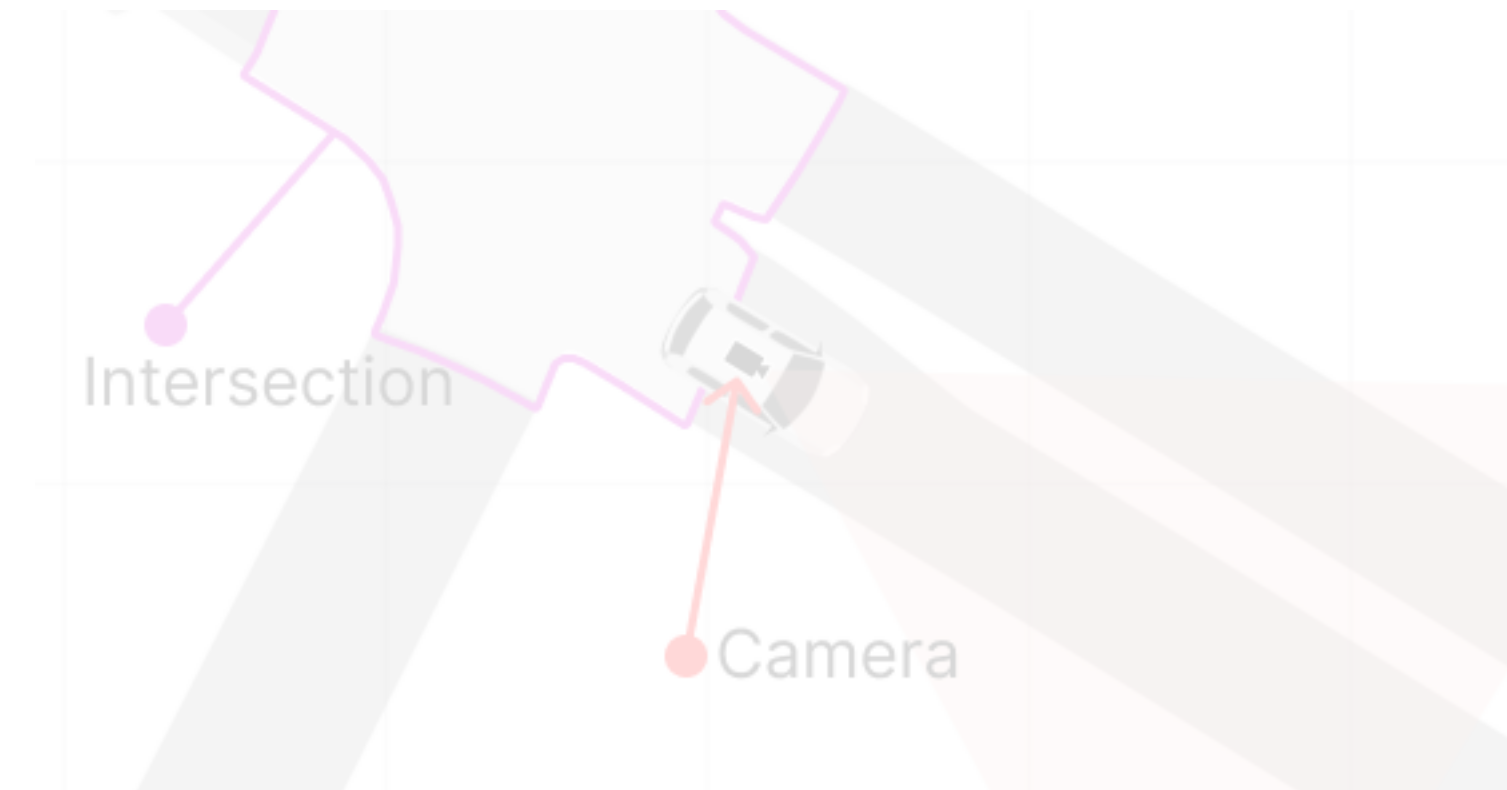
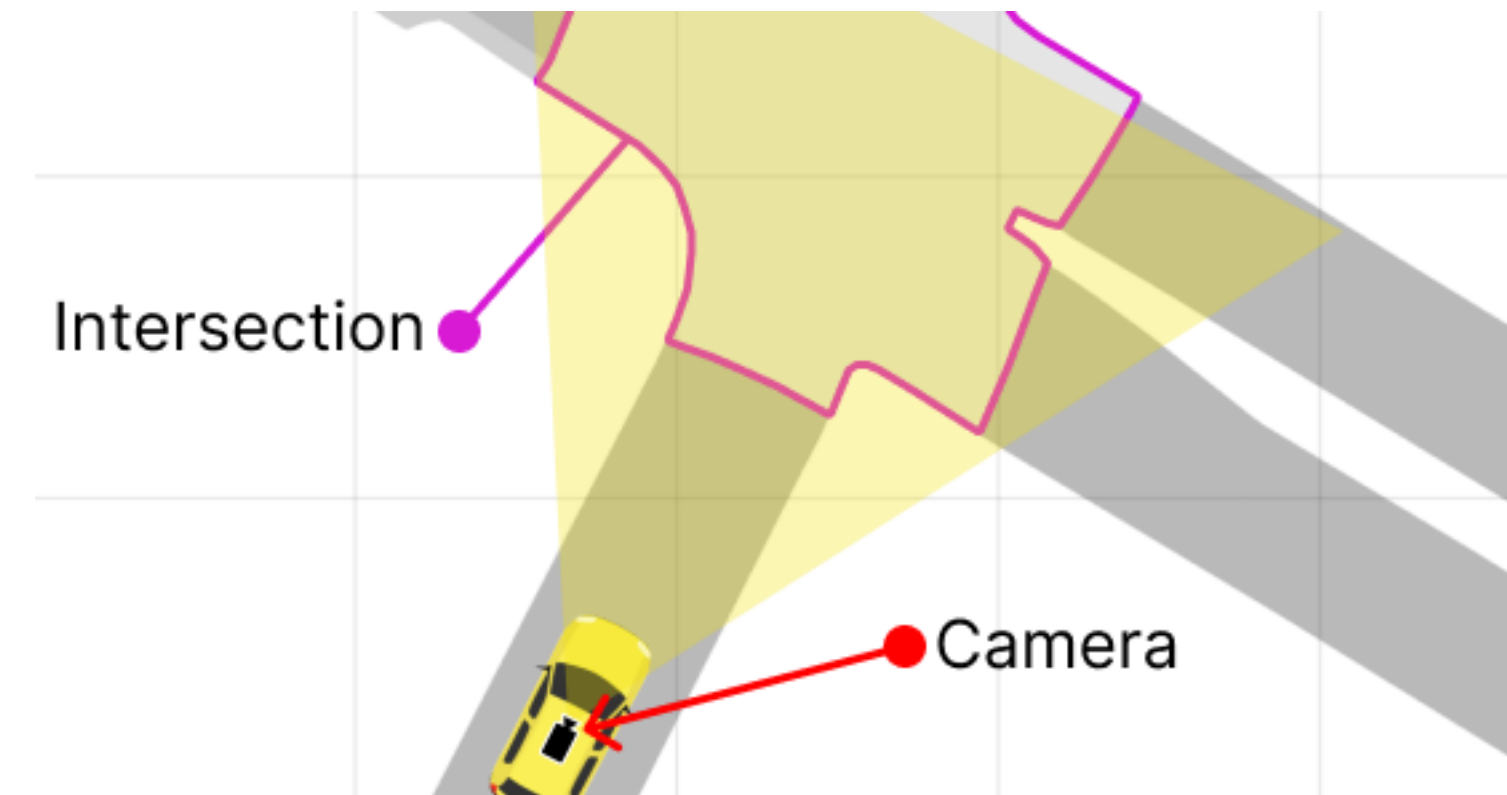
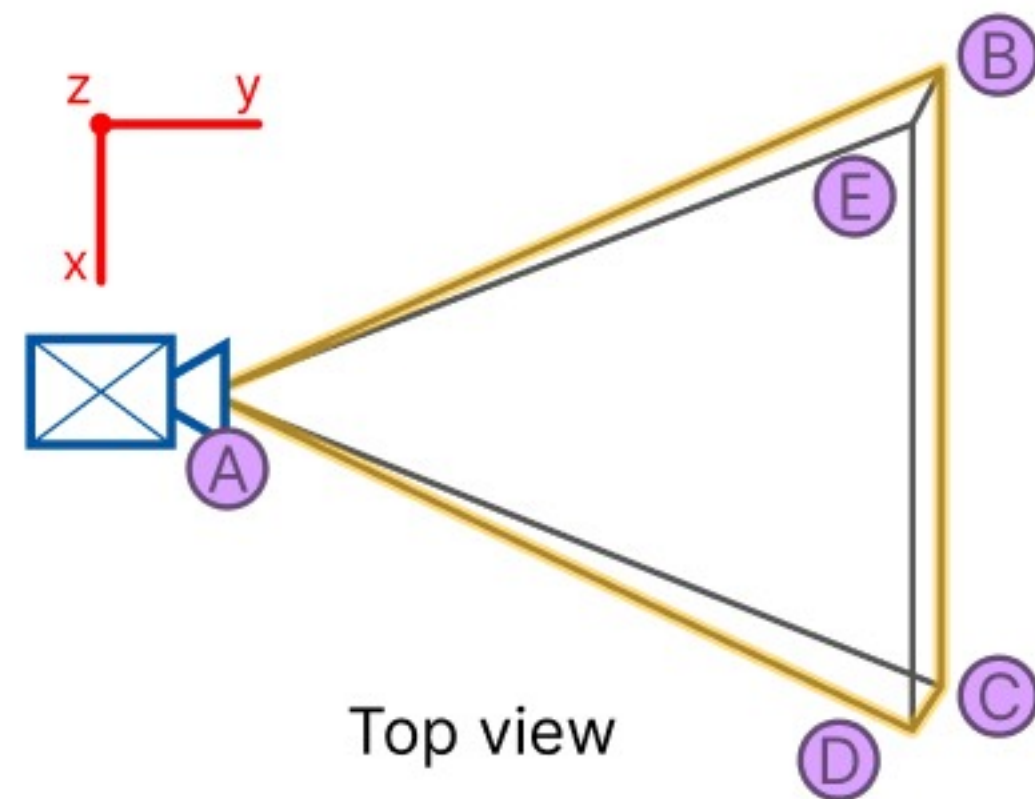
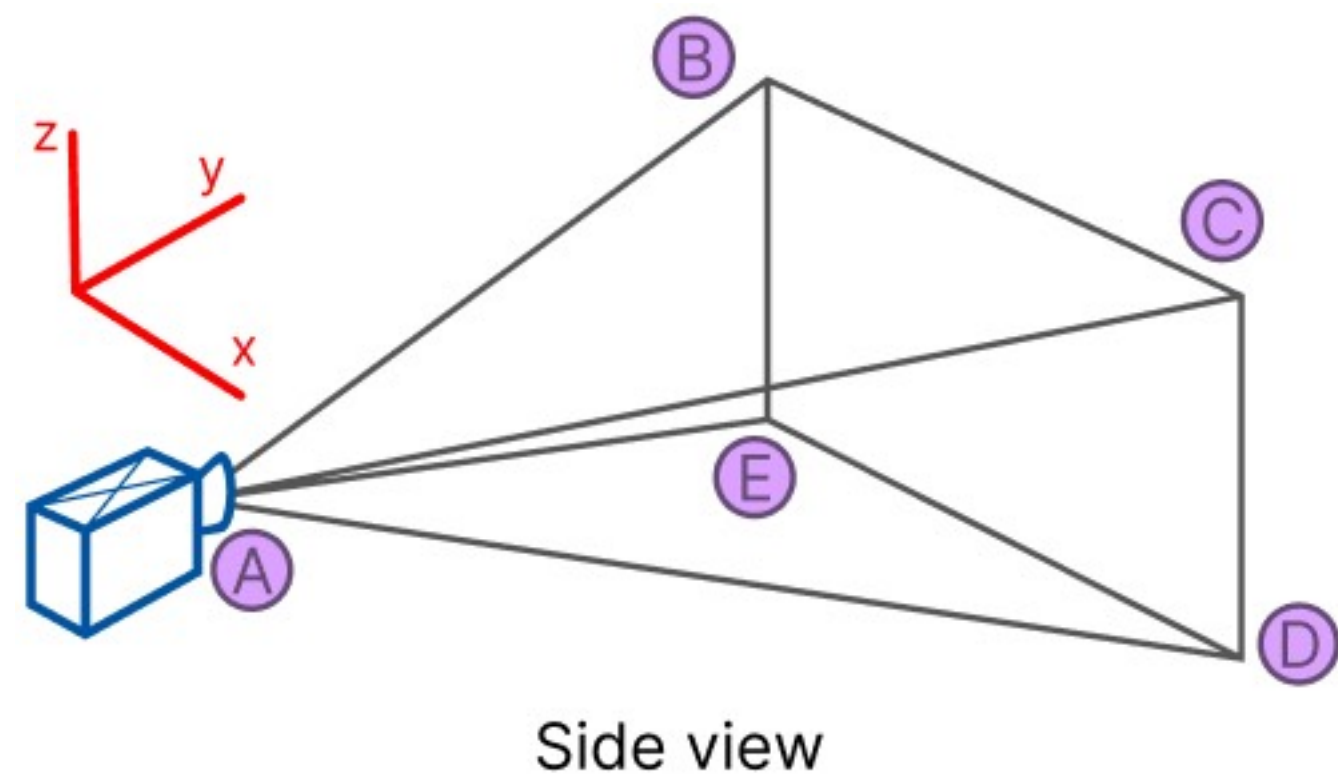
Optimization Techniques 01 Road Visibility Pruner

"Geographic constructs' visibility is a proxy for objects' visibility"



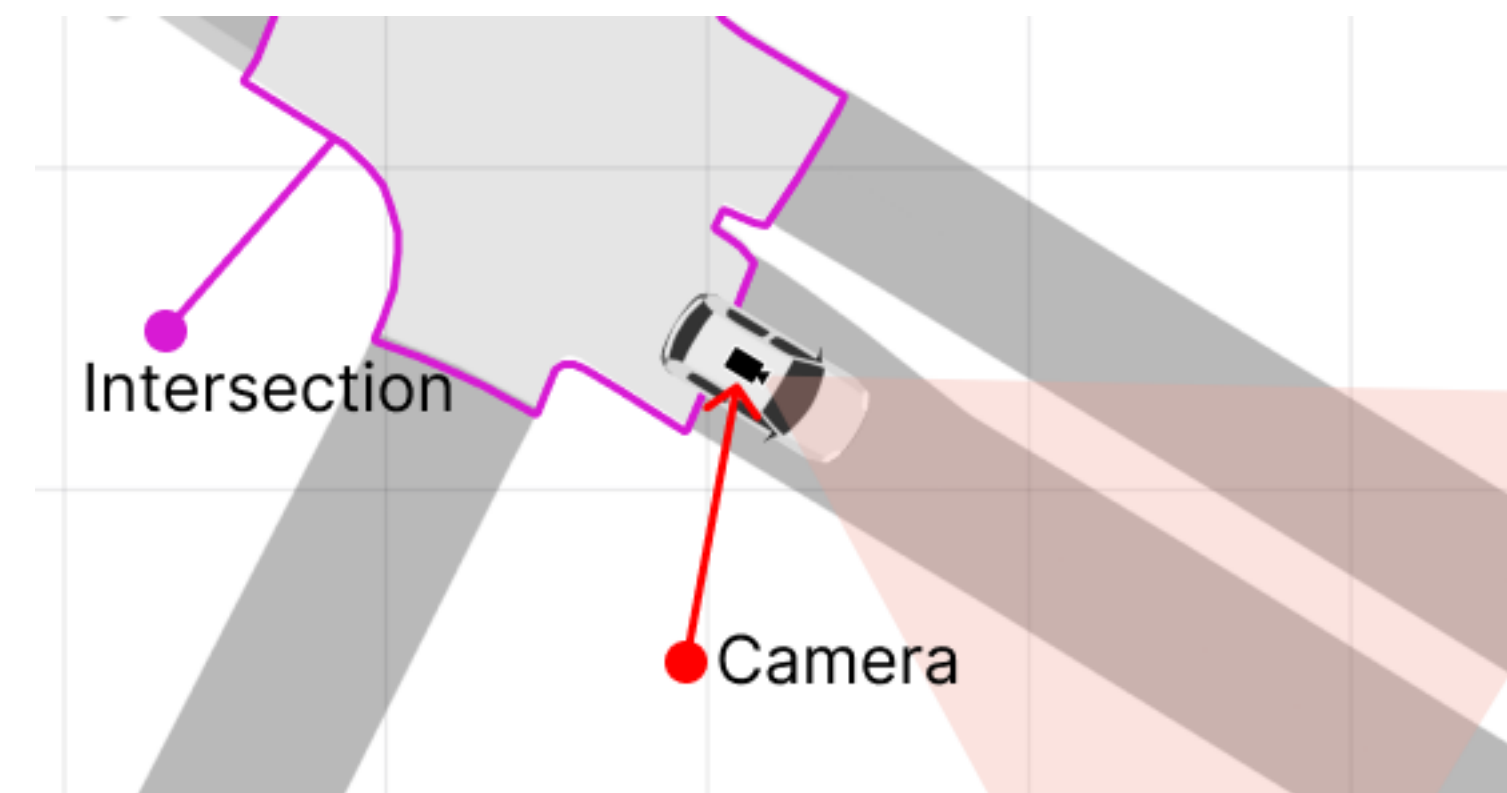
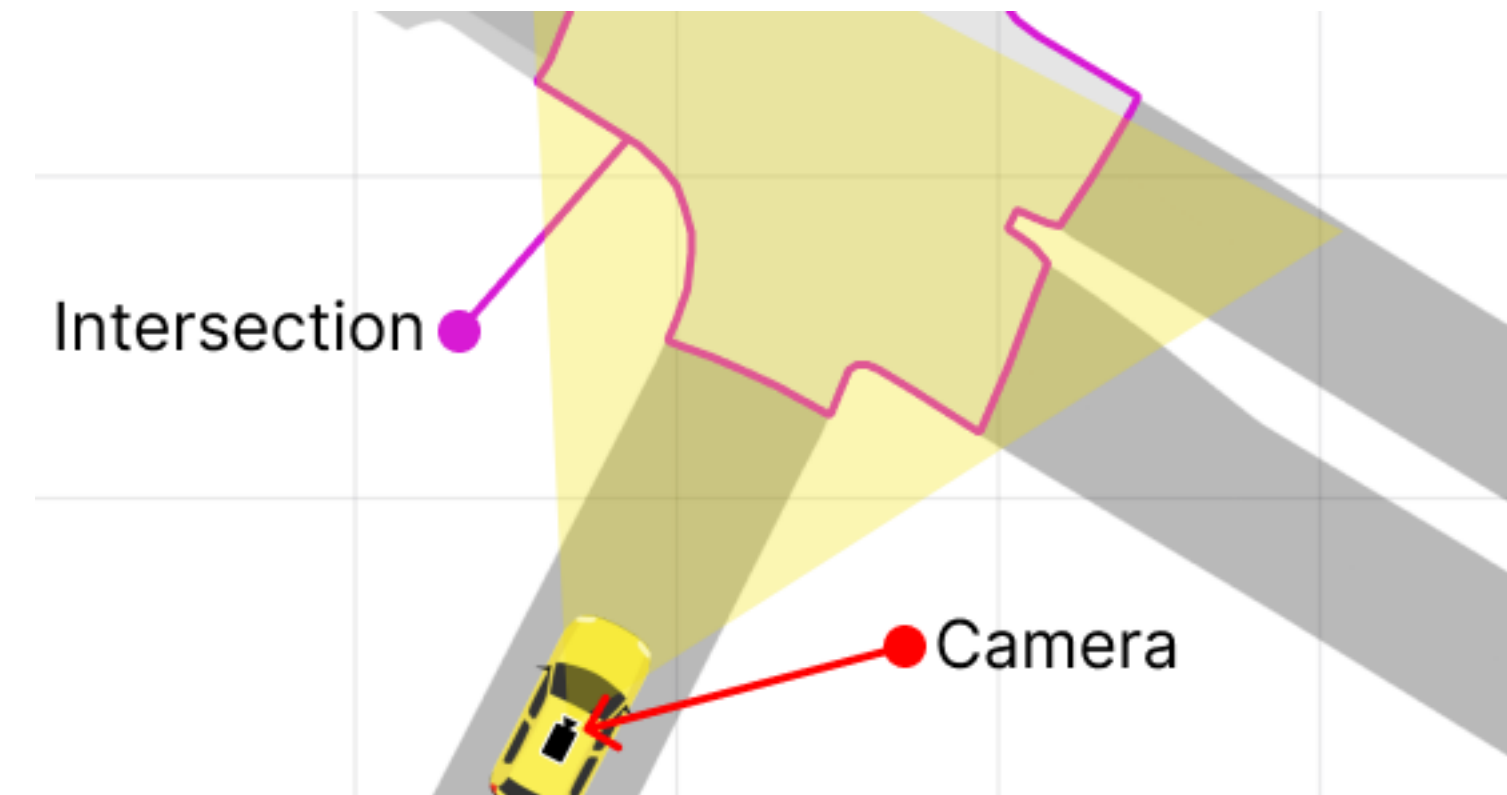
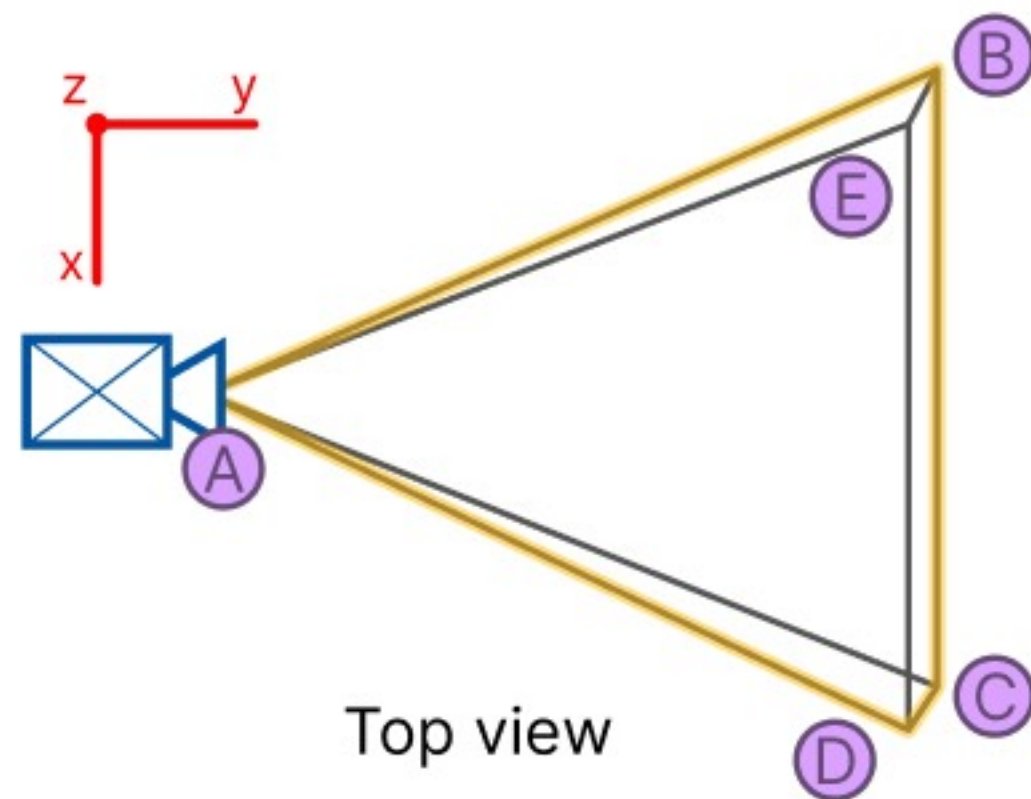
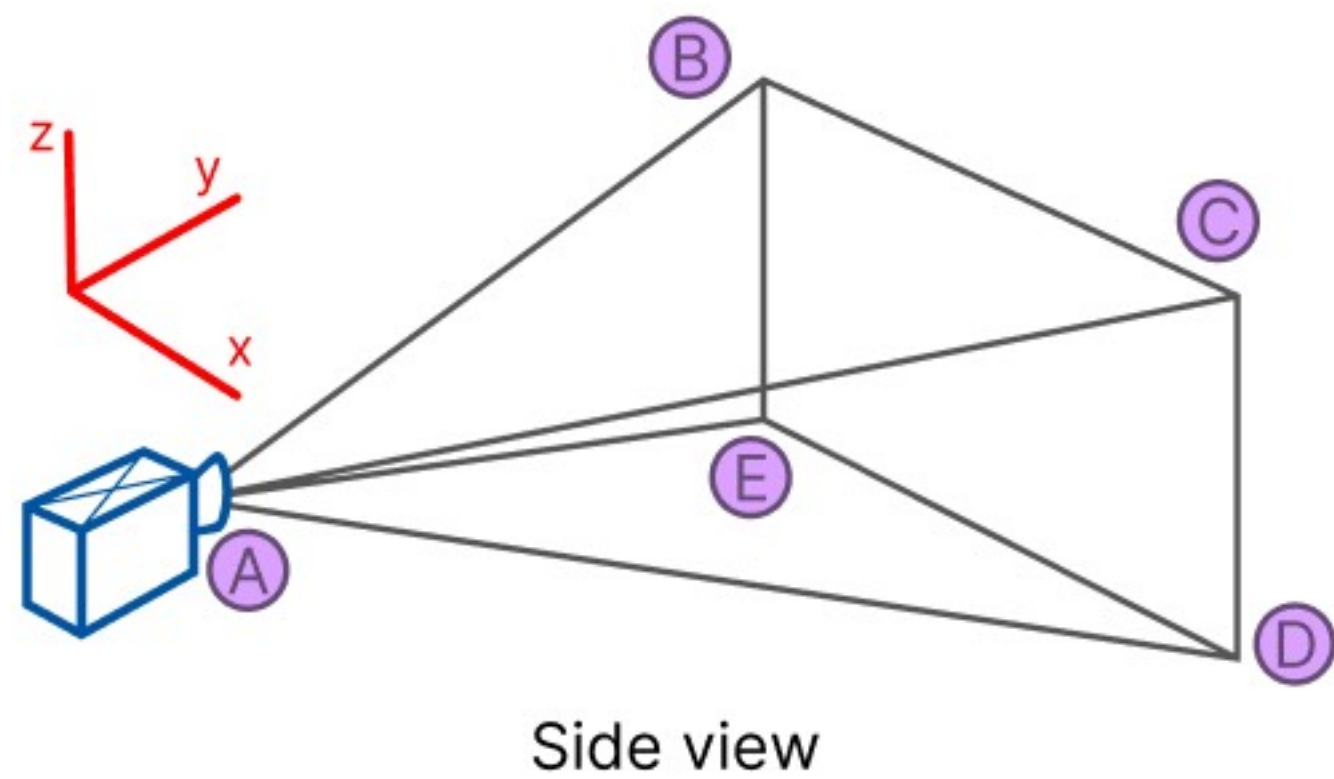
Optimization Techniques 01 Road Visibility Pruner

"Geographic constructs' visibility is a proxy for objects' visibility"



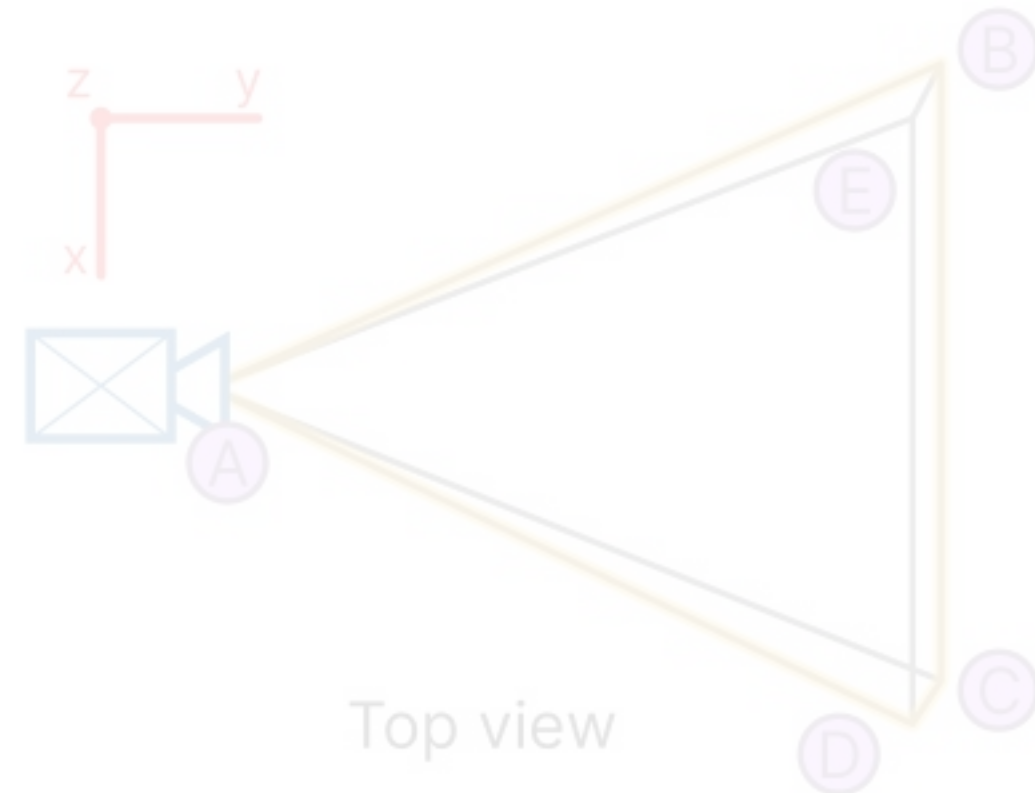
Optimization Techniques 01 Road Visibility Pruner

"Geographic constructs' visibility is a proxy for objects' visibility"



Optimization Techniques 01 Road Visibility Pruner

"Geographic constructs' visibility is a proxy for objects' visibility"



1.25x Speed up

Prunes out 21% of total frames



Optimization Techniques

03 3D Location Estimator (with geospatial metadata)

QUERY

```
(o.type=='car') | (o.type=='person')
```

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"

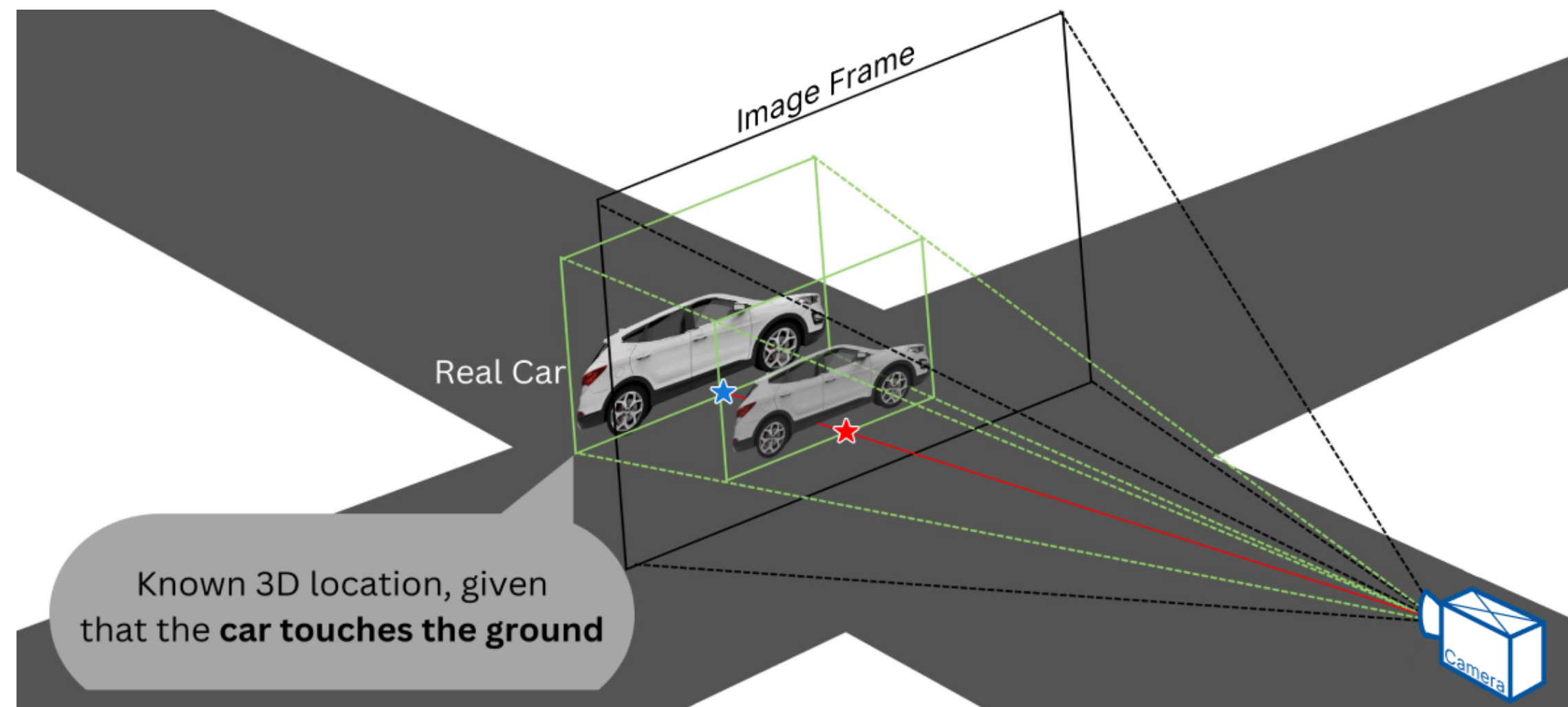
Optimization Techniques

03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"

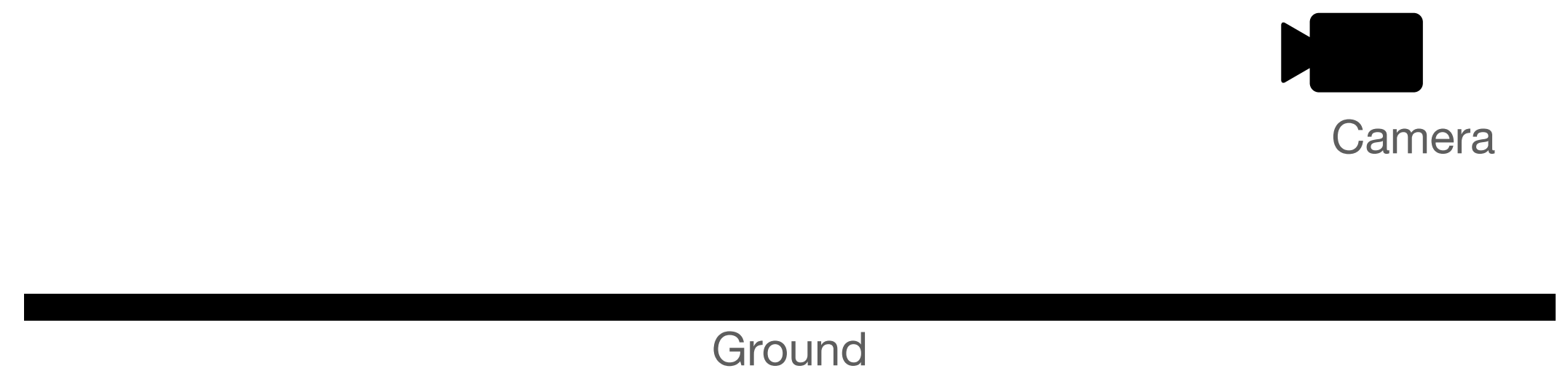
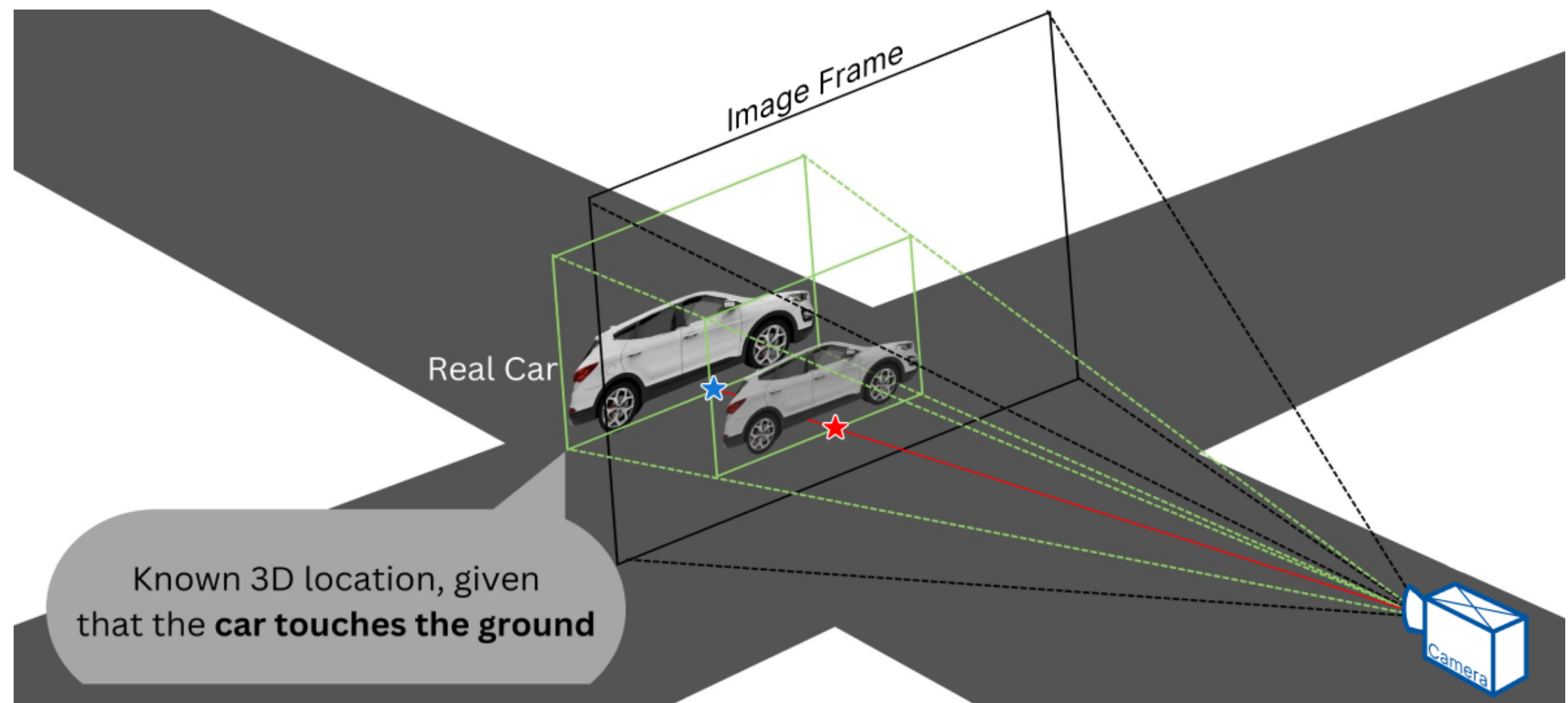
Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



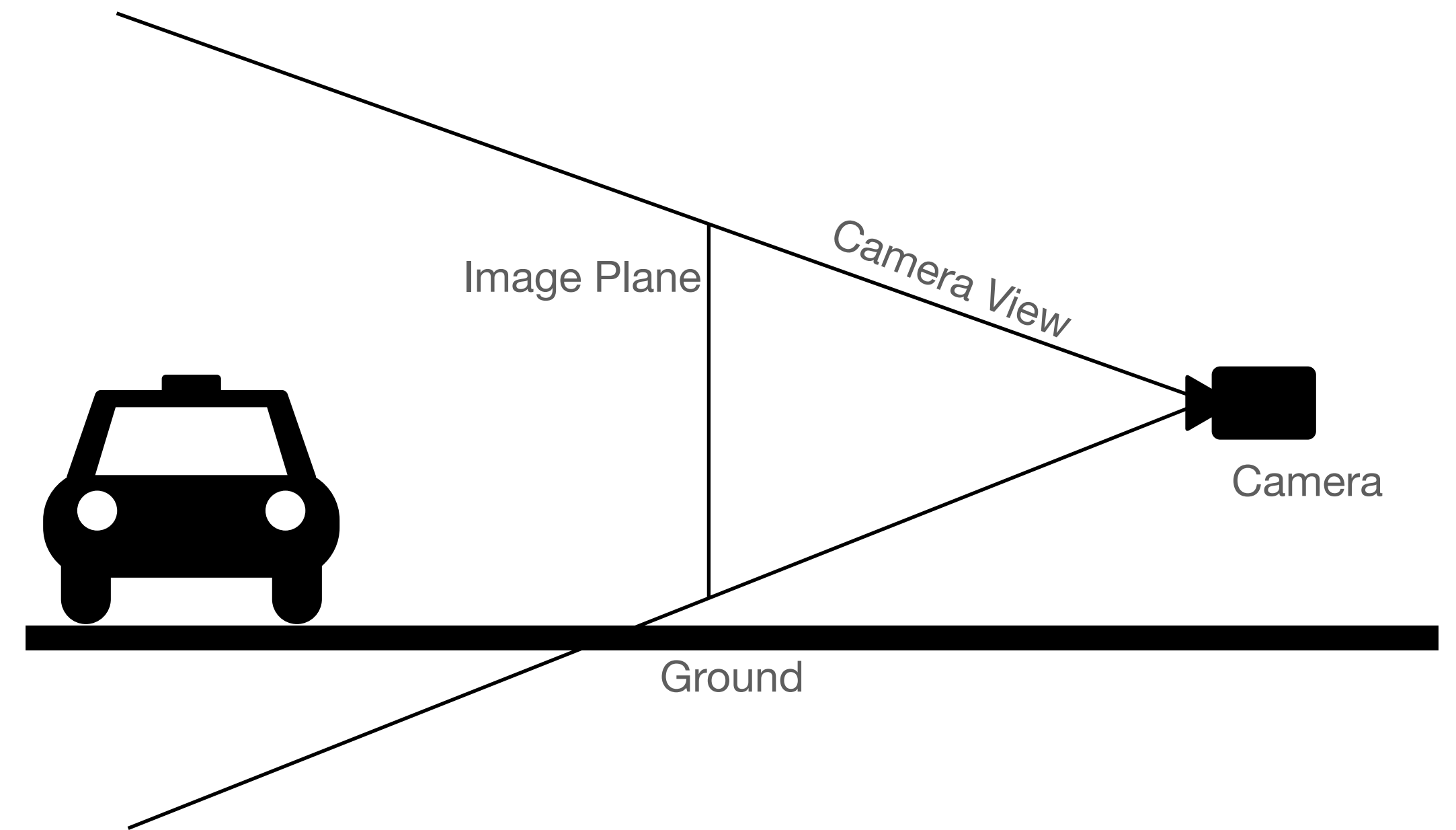
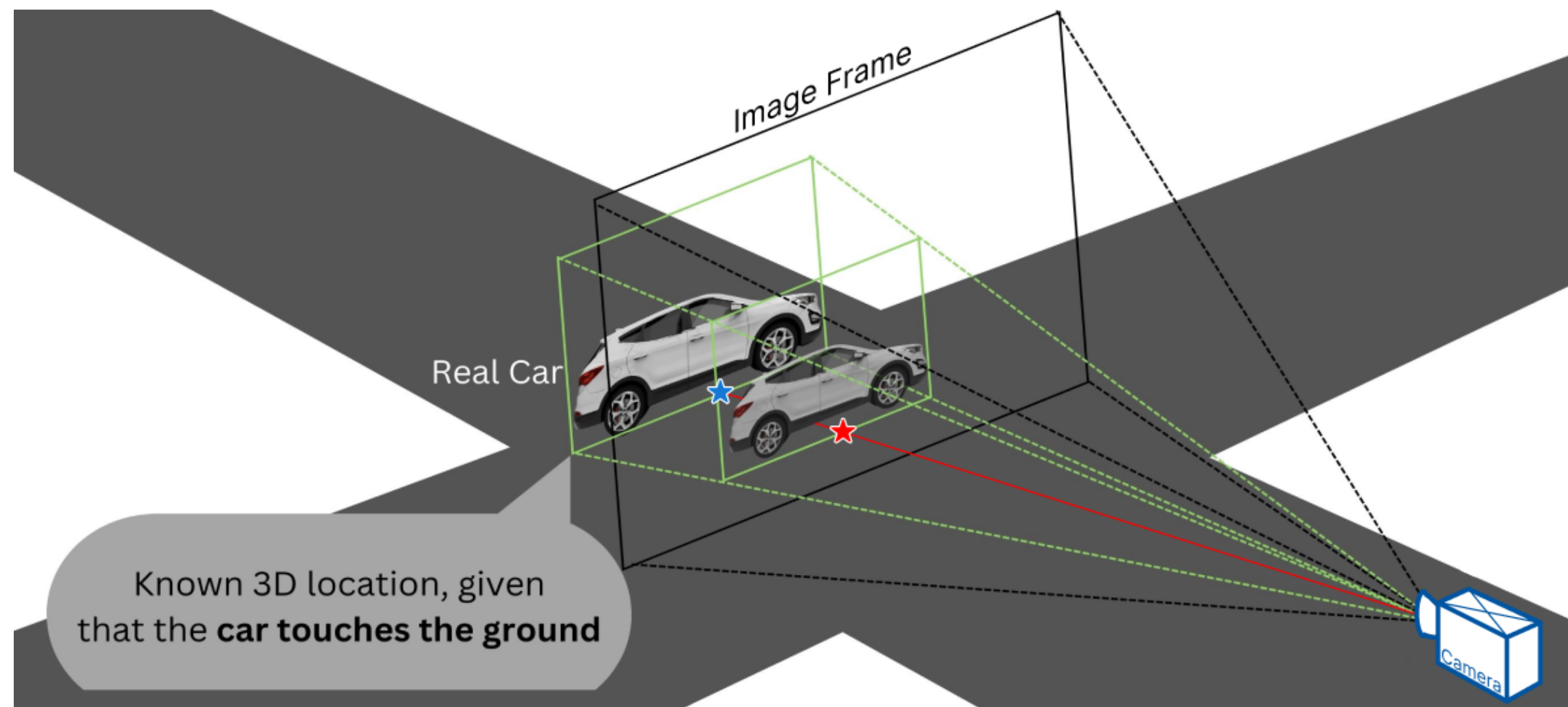
Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



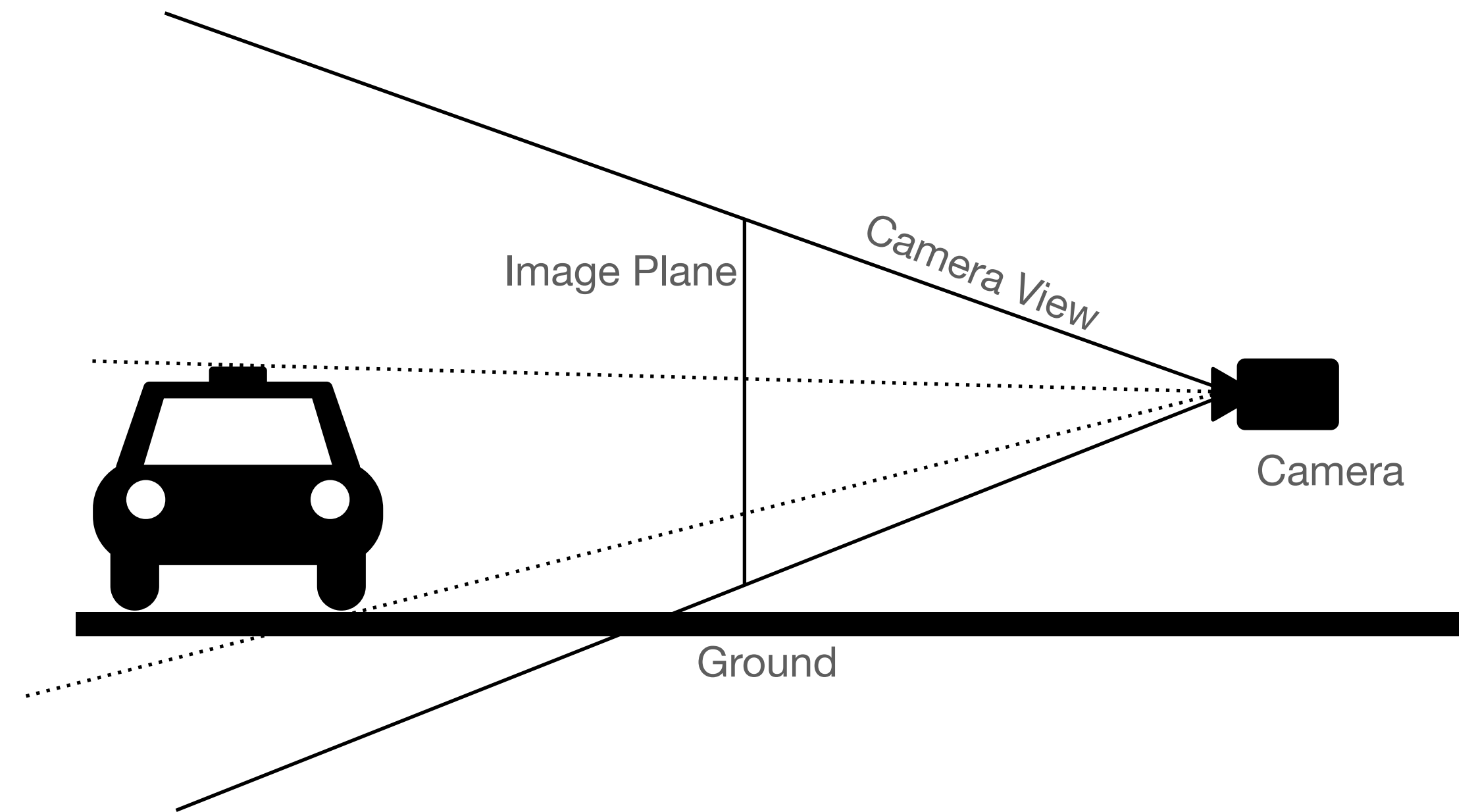
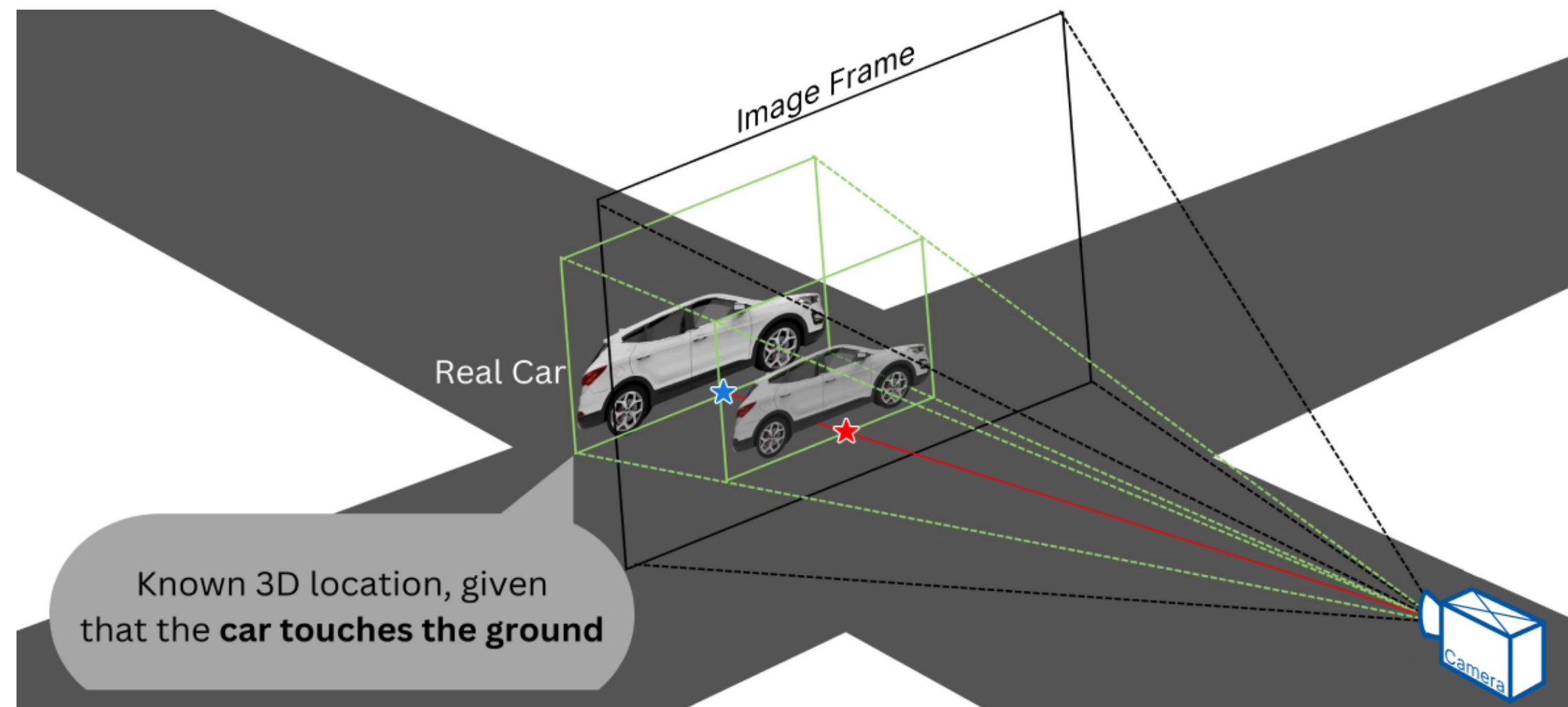
Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



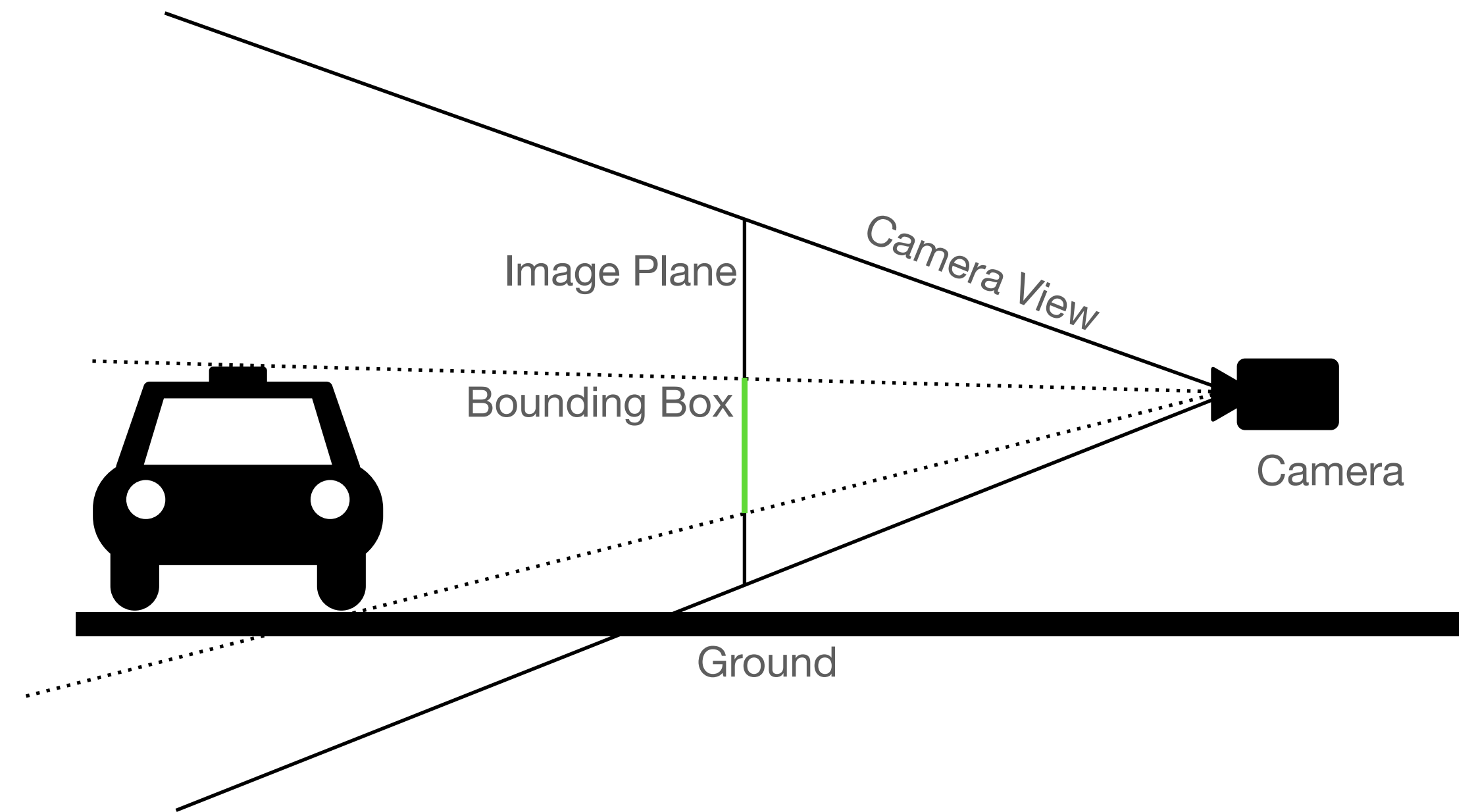
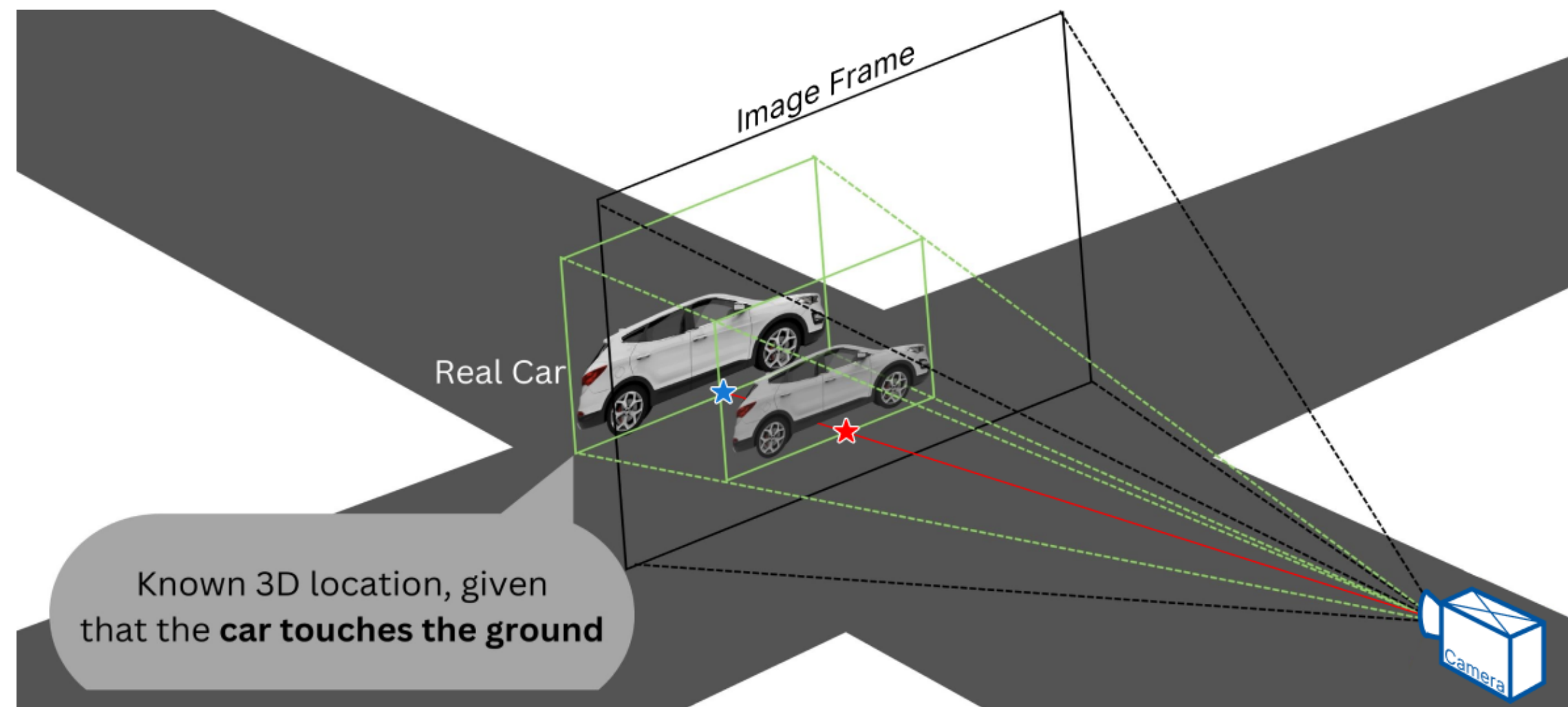
Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



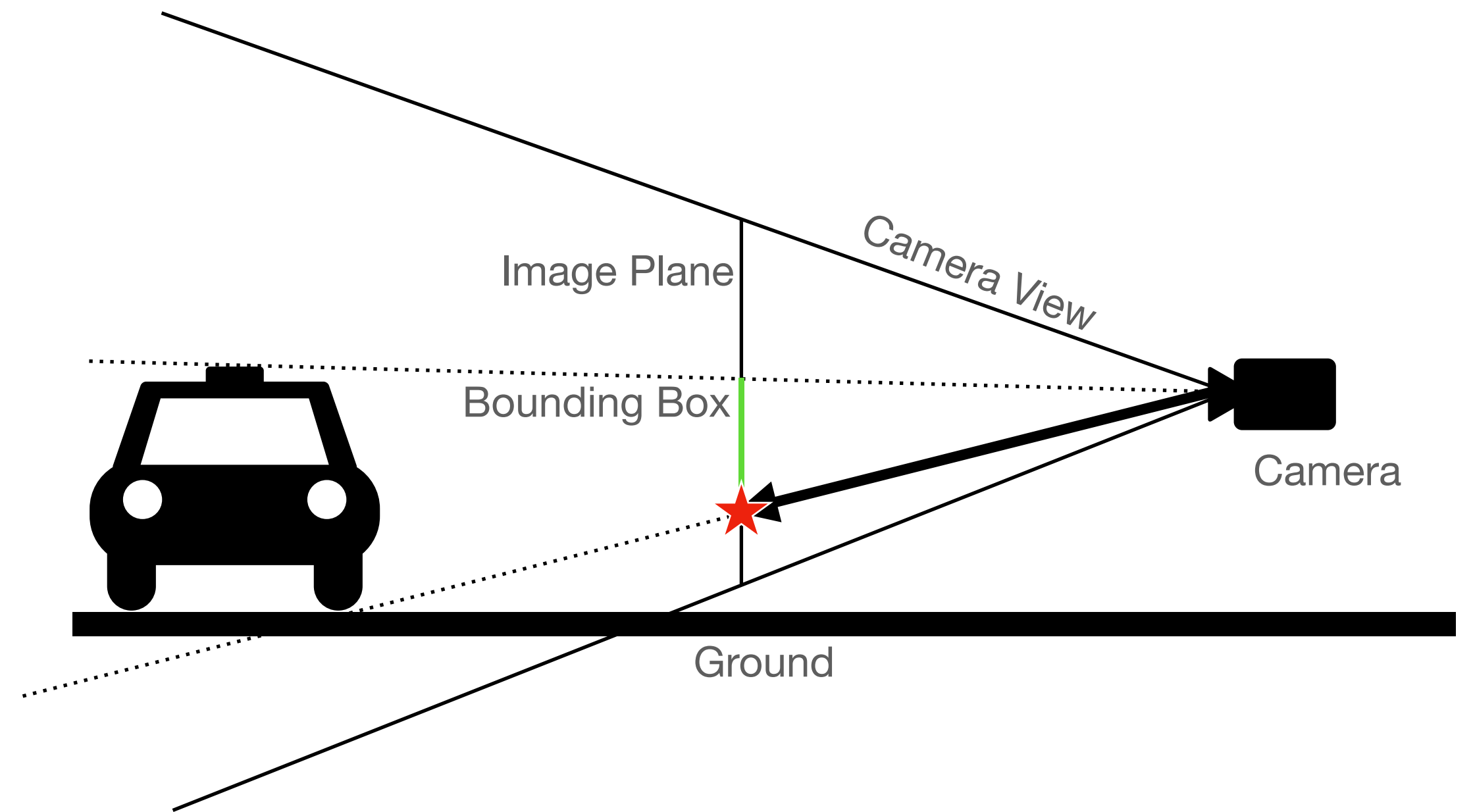
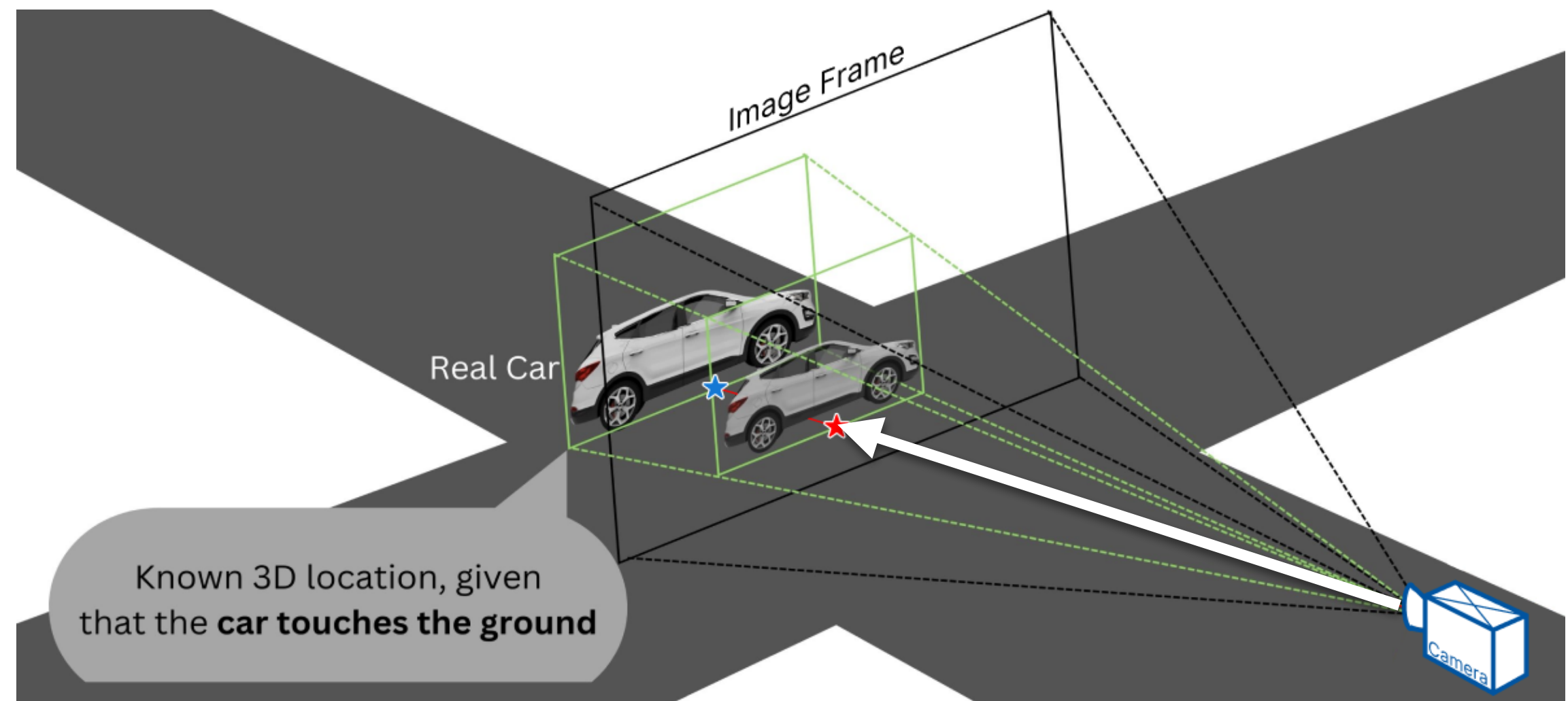
Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



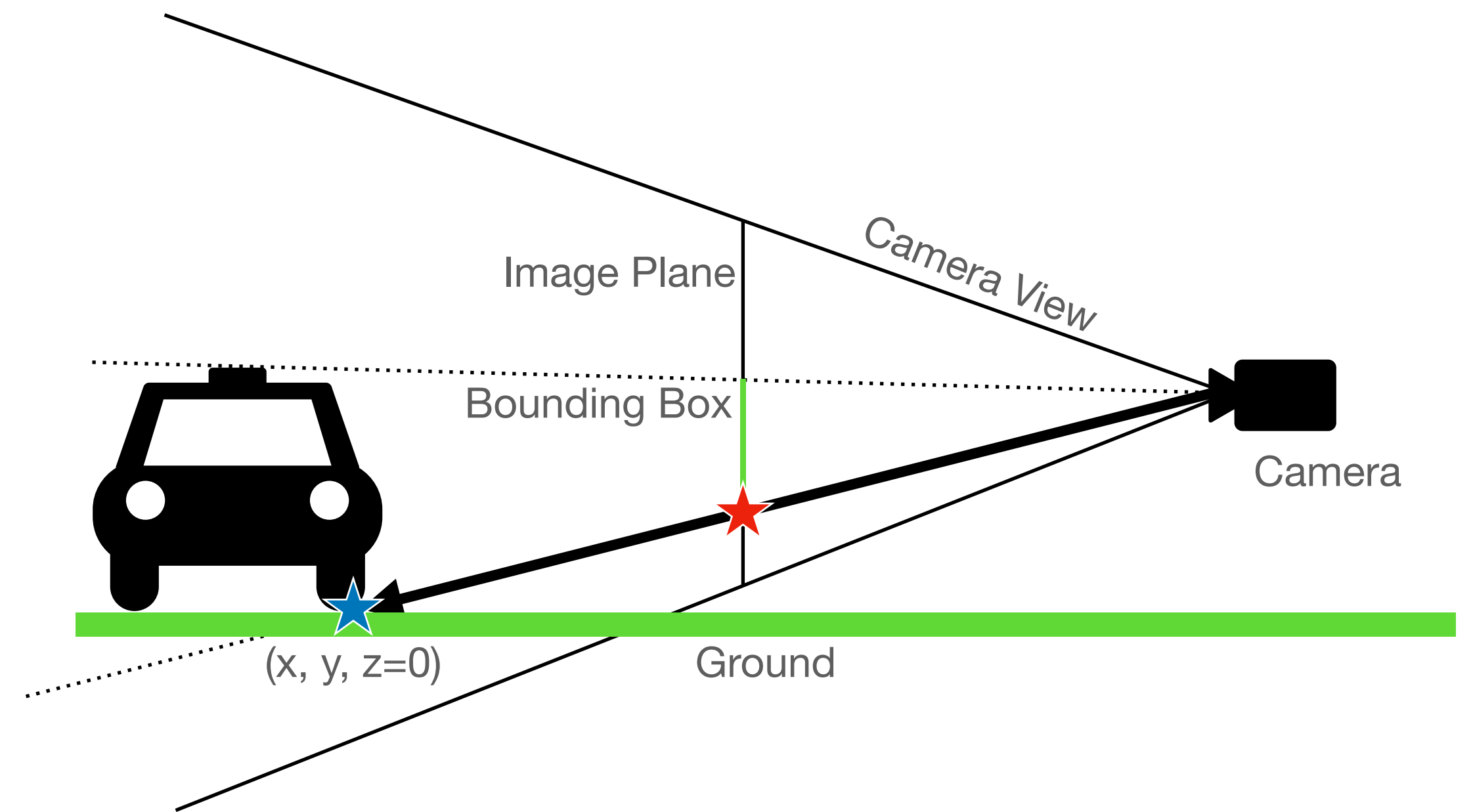
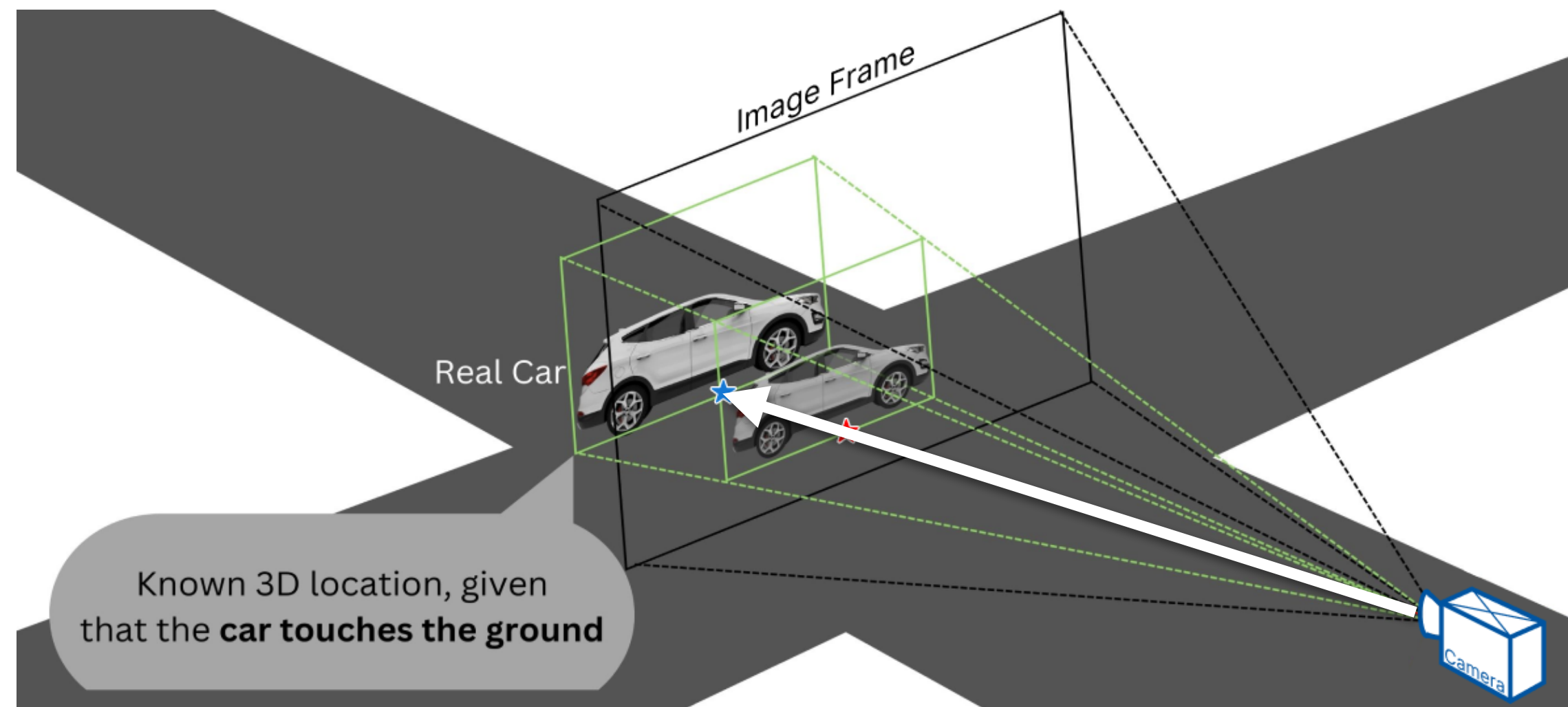
Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



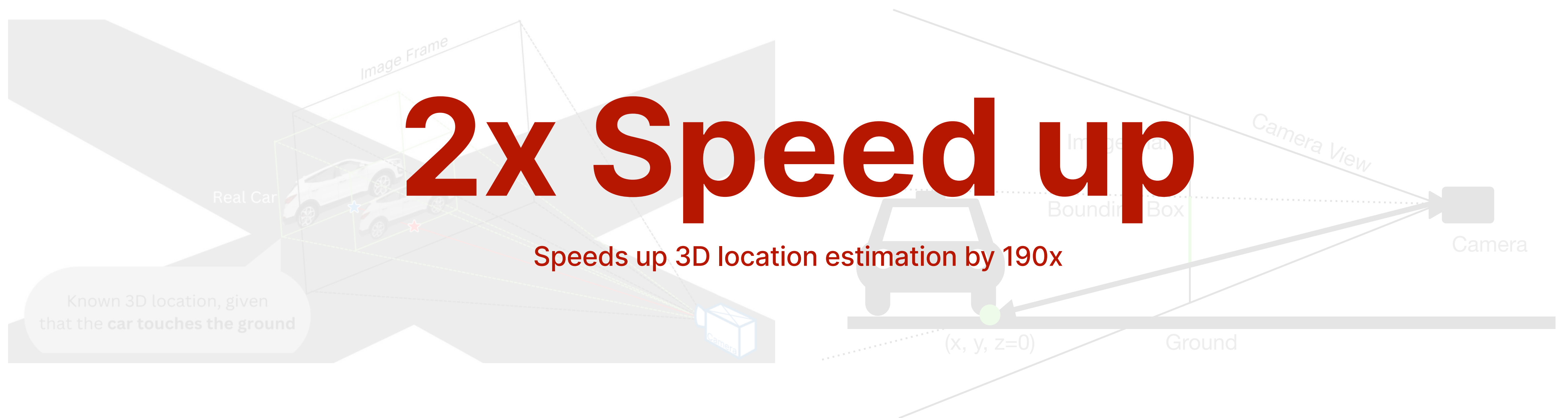
Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



Optimization Techniques 03 3D Location Estimator (with geospatial metadata)

"If objects of interest are on the ground, we can recover their 3D location based on their 2D pixels bounding box and camera parameters"



SPATIALYZE

#4 Evaluation

- * Experiment Setups
- * Results

Experiment Setups

VIVA

End-to-End video analytic system

VIVA: An End-to-End System for Interactive Video Analytics

Daniel Kang*
ddkang@cs.stanford.edu
Stanford University

Francisco Romero*
faromero@stanford.edu
Stanford University

Peter Bailis
pbailis@cs.stanford.edu
Stanford University

Christos Kozyrakis
christos@cs.stanford.edu
Stanford University

Matei Zaharia
matei@cs.stanford.edu
Stanford University and Databricks

ABSTRACT

The growth of video volumes and increased DNN capabilities has led to a growing desire for video analytics. In response, the data analytics community has proposed multiple systems that optimize specific query types (e.g., selection queries) or a particular step in query execution (e.g., video retrieval from storage). However, none of these systems provide *end-to-end, practical* video analytics for users to iteratively and interactively engage with queries, as is the case with analytics systems for structured data.

In response, we are building VIVA: an end-to-end system for interactive video analytics. VIVA contains five novel components. First, VIVA uses *relational hints*, which allow users to express relationships between columns that are difficult to automatically infer (e.g., mentions of a person in a transcript can be used as a proxy for the person appearing in the video). Second, VIVA introduces a *mixed-data query optimizer* that optimizes queries across both structured and unstructured data. Third, VIVA features an *embedding cache* that decides which results/embeddings to store for future queries. Finally, VIVA co-optimizes storage and query execution with its *video file manager* and *accelerator-based execution engine*. The former decides how to pre-fetch/manage video, while the latter selects and manages heterogeneous hardware backends spanning the growing number of DNN accelerators. We describe the challenges and design requirements for VIVA's development and outline ongoing and future work for realizing VIVA.

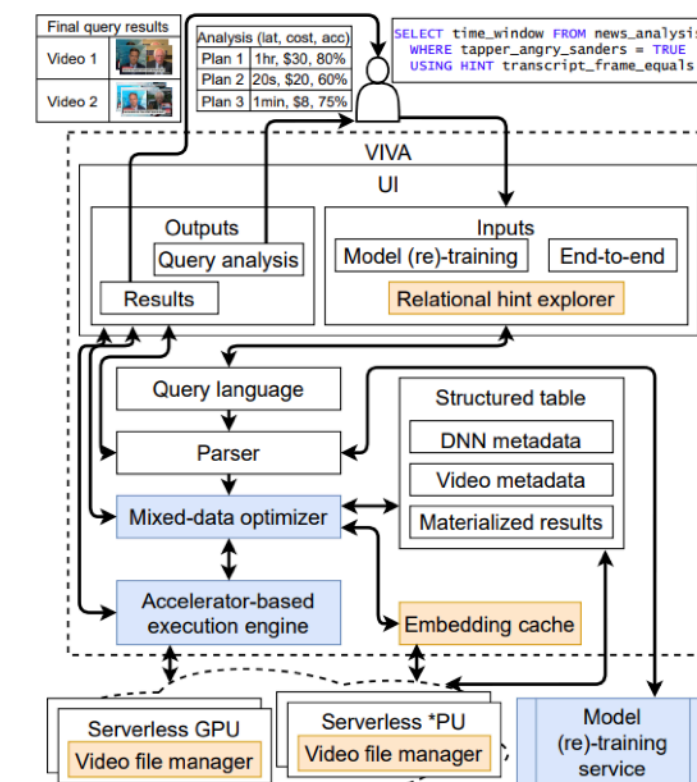


Figure 1: VIVA's architecture diagram. Blue components/interfaces are typically found in analytic systems for structured data that require rethinking. Orange components/interfaces are novel to VIVA.

SkyQuery

Geospatial video query processor

SkyQuery: An Aerial Drone Video Sensing Platform

Favyen Bastani
MIT CSAIL
US
favyen@csail.mit.edu

Songtao He
MIT CSAIL
US
songtao@csail.mit.edu

Ziwen Jiang
MIT CSAIL
US
ziwenj@csail.mit.edu

Osbert Bastani
University of Pennsylvania
US
obastani@seas.upenn.edu

Sam Madden
MIT CSAIL
US
madden@csail.mit.edu

Abstract

Video-based sensing from aerial drones, especially small multirotor drones, can provide rich data for numerous applications, including traffic analysis (computing traffic flow volumes), precision agriculture (periodically evaluating plant health), and wildlife population management (estimating population sizes). However, aerial drone video sensing applications must handle a surprisingly wide range of tasks: video frames must be aligned so that we can equate coordinates of objects that appear in different frames, video data must be analyzed to extract application-specific insights, and drone routes must be computed that maximize the value of newly captured video. To address these challenges, we built *SkyQuery*, a novel aerial drone video sensing platform that provides an expressive, high-level programming language to make it straightforward for users to develop complex long-running sensing applications. *SkyQuery* combines novel methods for fast video frame alignment and detection of small objects in top-down aerial drone video to efficiently execute applications with diverse video analysis workflows and data distributions, thereby allowing application developers to focus on the unique qualities of their particular application rather than general video processing, data analy-

ACM Reference Format:

Favyen Bastani, Songtao He, Ziwen Jiang, Osbert Bastani, and Sam Madden. 2021. *SkyQuery: An Aerial Drone Video Sensing Platform*. In *Proceedings of the 2021 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward'21)*, October 20–22, 2021, Chicago, IL, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3486607.3486750>

1 Introduction

When equipped with sensors, computation, and communication capabilities, small multirotor aerial drones are an ideal mobile sensing platform for many applications. Indeed, video sensing on aerial drones has been proposed for a number of tasks, including *earth science*, for monitoring coastal sea ice and habitat loss [24]; *disaster emergency response*, for tracking floods, fires, and people in search-and-rescue [8]; *automated mapping*, for measuring traffic incidents, congestion, and parking [17]; *civil infrastructure monitoring*, for finding defects in bridges, buildings, and pipelines [25]; and *precision agriculture*, for micro-climate monitoring of plant growth [4].

However, realizing applications that involve extracting insights from aerial drone video remains difficult; in fact, every

Experiment Setups

VIVA

End-to-End video analytic system

VIVA: An End-to-End System for Interactive Video Analytics

Daniel Kang*
ddkang@cs.stanford.edu
Stanford University

Francisco Romero*
faromero@stanford.edu
Stanford University

Peter Bailis
pbailis@cs.stanford.edu
Stanford University

Christos Kozyrakis
christos@cs.stanford.edu
Stanford University

Matei Zaharia
matei@cs.stanford.edu
Stanford University and Databricks

ABSTRACT

The growth of video volumes and increased DNN capabilities has led to a growing desire for video analytics. In response, the data analytics community has proposed multiple systems that optimize specific query types (e.g., selection queries) or a particular step in query execution (e.g., video retrieval from storage). However, none of these systems provide *end-to-end, practical* video analytics for users to iteratively and interactively engage with queries, as is the case with analytics systems for structured data.

In response, we are building VIVA: an end-to-end system for interactive video analytics. VIVA contains five novel components. First, VIVA uses *relational hints*, which allow users to express relationships between columns that are difficult to automatically infer (e.g., mentions of a person in a transcript can be used as a proxy for the person appearing in the video). Second, VIVA introduces a *mixed-data query optimizer* that optimizes queries across both structured and unstructured data. Third, VIVA features an *embedding cache* that decides which results/embeddings to store for future queries. Finally, VIVA co-optimizes storage and query execution with its *video file manager* and *accelerator-based execution engine*. The former decides how to pre-fetch/manage video, while the latter selects and manages heterogeneous hardware backends spanning the growing number of DNN accelerators. We describe the challenges and design requirements for VIVA's development and outline ongoing and future work for realizing VIVA.

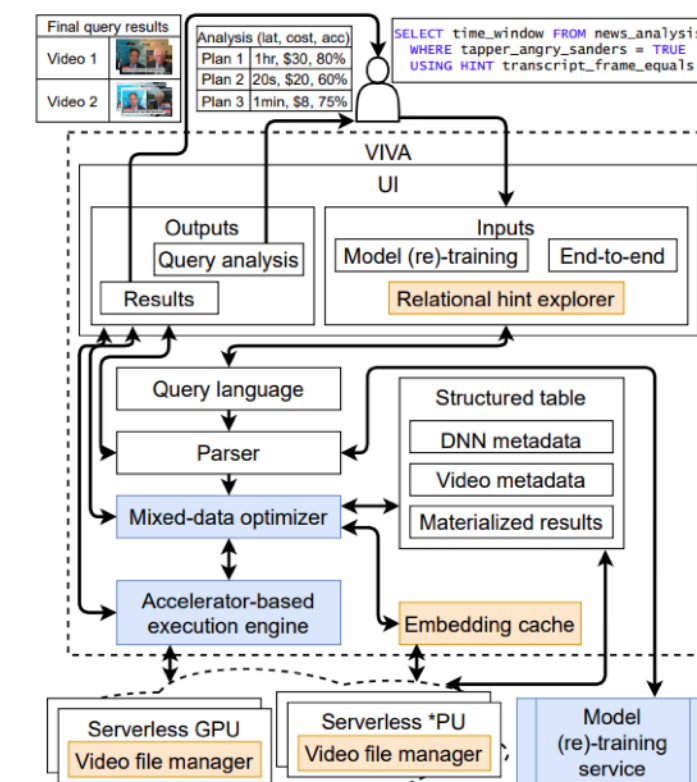


Figure 1: VIVA's architecture diagram. Blue components/interfaces are typically found in analytic systems for structured data that require rethinking. Orange components/interfaces are novel to VIVA.

SkyQuery

Geospatial video query processor

SkyQuery: An Aerial Drone Video Sensing Platform

Favyen Bastani
MIT CSAIL
US
favyen@csail.mit.edu

Songtao He
MIT CSAIL
US
songtao@csail.mit.edu

Ziwen Jiang
MIT CSAIL
US
ziwenj@csail.mit.edu

Osbert Bastani
University of Pennsylvania
US
obastani@seas.upenn.edu

Sam Madden
MIT CSAIL
US
madden@csail.mit.edu

Abstract

Video-based sensing from aerial drones, especially small multirotor drones, can provide rich data for numerous applications, including traffic analysis (computing traffic flow volumes), precision agriculture (periodically evaluating plant health), and wildlife population management (estimating population sizes). However, aerial drone video sensing applications must handle a surprisingly wide range of tasks: video frames must be aligned so that we can equate coordinates of objects that appear in different frames, video data must be analyzed to extract application-specific insights, and drone routes must be computed that maximize the value of newly captured video. To address these challenges, we built *SkyQuery*, a novel aerial drone video sensing platform that provides an expressive, high-level programming language to make it straightforward for users to develop complex long-running sensing applications. *SkyQuery* combines novel methods for fast video frame alignment and detection of small objects in top-down aerial drone video to efficiently execute applications with diverse video analysis workflows and data distributions, thereby allowing application developers to focus on the unique qualities of their particular application rather than general video processing, data analy-

ACM Reference Format:

Favyen Bastani, Songtao He, Ziwen Jiang, Osbert Bastani, and Sam Madden. 2021. *SkyQuery: An Aerial Drone Video Sensing Platform*. In *Proceedings of the 2021 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward'21)*, October 20–22, 2021, Chicago, IL, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3486607.3486750>

1 Introduction

When equipped with sensors, computation, and communication capabilities, small multirotor aerial drones are an ideal mobile sensing platform for many applications. Indeed, video sensing on aerial drones has been proposed for a number of tasks, including *earth science*, for monitoring coastal sea ice and habitat loss [24]; *disaster emergency response*, for tracking floods, fires, and people in search-and-rescue [8]; *automated mapping*, for measuring traffic incidents, congestion, and parking [17]; *civil infrastructure monitoring*, for finding defects in bridges, buildings, and pipelines [25]; and *precision agriculture*, for micro-climate monitoring of plant growth [4].

However, realizing applications that involve extracting insights from aerial drone video remains difficult; in fact, every

Experiment Setups

VIVA

End-to-End video analytic system

VIVA: An End-to-End System for Interactive Video Analytics

Daniel Kang*
ddkang@cs.stanford.edu
Stanford University

Francisco Romero*
faromero@stanford.edu
Stanford University

Peter Bailis
pbailis@cs.stanford.edu
Stanford University

Christos Kozyrakis
christos@cs.stanford.edu
Stanford University

Matei Zaharia
matei@cs.stanford.edu
Stanford University and Databricks

ABSTRACT

The growth of video volumes and increased DNN capabilities has led to a growing desire for video analytics. In response, the data analytics community has proposed multiple systems that optimize specific query types (e.g., selection queries) or a particular step in query execution (e.g., video retrieval from storage). However, none of these systems provide *end-to-end, practical* video analytics for users to iteratively and interactively engage with queries, as is the case with analytics systems for structured data.

In response, we are building VIVA: an end-to-end system for interactive video analytics. VIVA contains five novel components. First, VIVA uses *relational hints*, which allow users to express relationships between columns that are difficult to automatically infer (e.g., mentions of a person in a transcript can be used as a proxy for the person appearing in the video). Second, VIVA introduces a *mixed-data query optimizer* that optimizes queries across both structured and unstructured data. Third, VIVA features an *embedding cache* that decides which results/embeddings to store for future queries. Finally, VIVA co-optimizes storage and query execution with its *video file manager* and *accelerator-based execution engine*. The former decides how to pre-fetch/manage video, while the latter selects and manages heterogeneous hardware backends spanning the growing number of DNN accelerators. We describe the challenges and design requirements for VIVA's development and outline ongoing and future work for realizing VIVA.

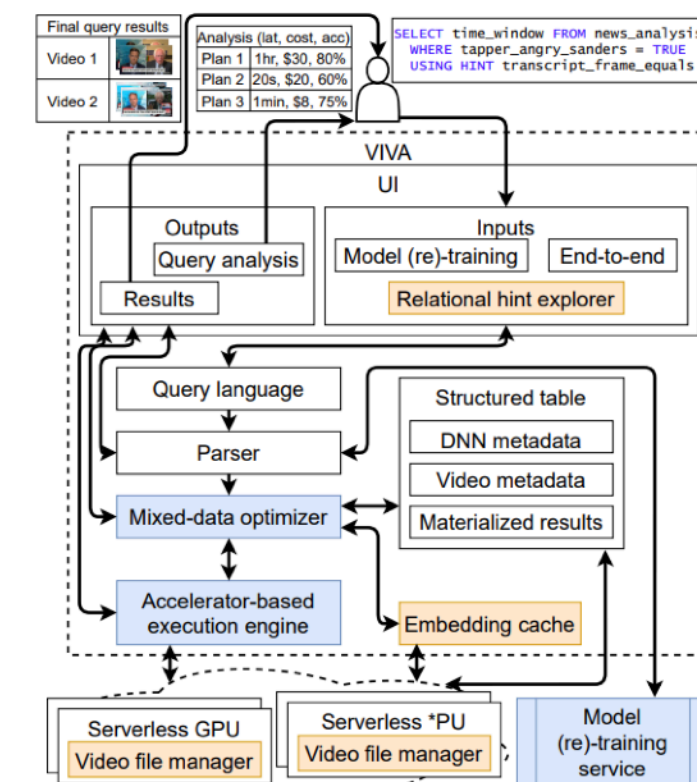


Figure 1: VIVA's architecture diagram. Blue components/interfaces are typically found in analytic systems for structured data that require rethinking. Orange components/interfaces are novel to VIVA.

SkyQuery

Geospatial video query processor

SkyQuery: An Aerial Drone Video Sensing Platform

Favyen Bastani
MIT CSAIL
US
favyen@csail.mit.edu

Songtao He
MIT CSAIL
US
songtao@csail.mit.edu

Ziwen Jiang
MIT CSAIL
US
ziwenj@csail.mit.edu

Osbert Bastani
University of Pennsylvania
US
obastani@seas.upenn.edu

Sam Madden
MIT CSAIL
US
madden@csail.mit.edu

Abstract

Video-based sensing from aerial drones, especially small multirotor drones, can provide rich data for numerous applications, including traffic analysis (computing traffic flow volumes), precision agriculture (periodically evaluating plant health), and wildlife population management (estimating population sizes). However, aerial drone video sensing applications must handle a surprisingly wide range of tasks: video frames must be aligned so that we can equate coordinates of objects that appear in different frames, video data must be analyzed to extract application-specific insights, and drone routes must be computed that maximize the value of newly captured video. To address these challenges, we built *SkyQuery*, a novel aerial drone video sensing platform that provides an expressive, high-level programming language to make it straightforward for users to develop complex long-running sensing applications. *SkyQuery* combines novel methods for fast video frame alignment and detection of small objects in top-down aerial drone video to efficiently execute applications with diverse video analysis workflows and data distributions, thereby allowing application developers to focus on the unique qualities of their particular application rather than general video processing, data analy-

ACM Reference Format:

Favyen Bastani, Songtao He, Ziwen Jiang, Osbert Bastani, and Sam Madden. 2021. *SkyQuery: An Aerial Drone Video Sensing Platform*. In *Proceedings of the 2021 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward'21)*, October 20–22, 2021, Chicago, IL, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3486607.3486750>

1 Introduction

When equipped with sensors, computation, and communication capabilities, small multirotor aerial drones are an ideal mobile sensing platform for many applications. Indeed, video sensing on aerial drones has been proposed for a number of tasks, including *earth science*, for monitoring coastal sea ice and habitat loss [24]; *disaster emergency response*, for tracking floods, fires, and people in search-and-rescue [8]; *automated mapping*, for measuring traffic incidents, congestion, and parking [17]; *civil infrastructure monitoring*, for finding defects in bridges, buildings, and pipelines [25]; and *precision agriculture*, for micro-climate monitoring of plant growth [4].

However, realizing applications that involve extracting insights from aerial drone video remains difficult; in fact, every

Experiment Setups

Ablation Compare with no Optimization

Q1 A pedestrian at an intersection moving perpendicularly to the camera

```
p = w.object(); c = w.camera(); i = w.geoConstruct('intersection')
w.filter(p.type=='person' & perpendicular(p, c)
        & contains(i, p))
```

SkyQuery Geospatial video query processor

Q3 A car stopped in a cycling lane.

```
c = w.camera(); b = w.geoConstruct('bike-lane')
w.filter(p.type=='car' & contains(b, c))
```

VIVA End-to-End video analytic system

Q2 A car turning left with a pedestrian at an intersection.

```
c = w.object(); p = w.object(); i = w.geoConstruct('intersection')
w.filter(c.type=='car' & p.type='people' & turnLeft(c)
        & contains(i, [c, p]))
```

Experiment Setups

Ablation Compare with no Optimization

Q1 A pedestrian at an intersection moving perpendicularly to the camera

```
p = w.object(); c = w.camera(); i = w.geoConstruct('intersection')
w.filter(p.type=='person' & perpendicular(p, c)
        & contains(i, p))
```

SkyQuery Geospatial video query processor

Q3 A car stopped in a cycling lane.

```
c = w.camera(); b = w.geoConstruct('bike-lane')
w.filter(p.type=='car' & contains(b, c))
```

VIVA End-to-End video analytic system

Q2 A car turning left with a pedestrian at an intersection.

```
c = w.object(); p = w.object(); i = w.geoConstruct('intersection')
w.filter(c.type=='car' & p.type='people' & turnLeft(c)
        & contains(i, [c, p]))
```


Experiment Setups

Ablation Compare with no Optimization

Q1 A pedestrian at an intersection moving perpendicularly to the camera

```
p = w.object(); c = w.camera(); i = w.geoConstruct('intersection')
w.filter(p.type=='person' & perpendicular(p, c)
        & contains(i, p))
```

SkyQuery Geospatial video query processor

Q3 A car stopped in a cycling lane.

```
c = w.camera(); b = w.geoConstruct('bike-lane')
w.filter(p.type=='car' & contains(b, c))
```

VIVA End-to-End video analytic system

Q2 A car turning left with a pedestrian at an intersection.

```
c = w.object(); p = w.object(); i = w.geoConstruct('intersection')
w.filter(c.type=='car' & p.type='people' & turnLeft(c)
        & contains(i, [c, p]))
```

Experiment Setups

Ablation Compare with no Optimization

Q1 A pedestrian at an intersection moving perpendicularly to the camera

```
p = w.object(); c = w.camera(); i = w.geoConstruct('intersection')
w.filter(p.type=='person' & perpendicular(p, c)
        & contains(i, p))
```

SkyQuery Geospatial video query processor

Q3 A car stopped in a cycling lane.

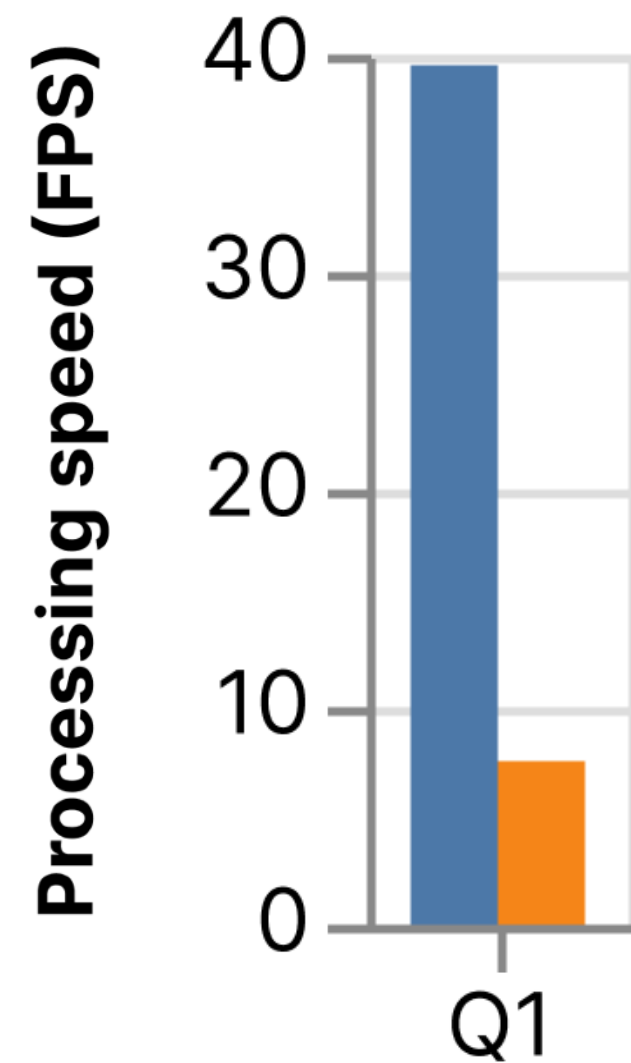
```
c = w.camera(); b = w.geoConstruct('bike-lane')
w.filter(p.type=='car' & contains(b, c))
```

VIVA End-to-End video analytic system

Q2 A car turning left with a pedestrian at an intersection.

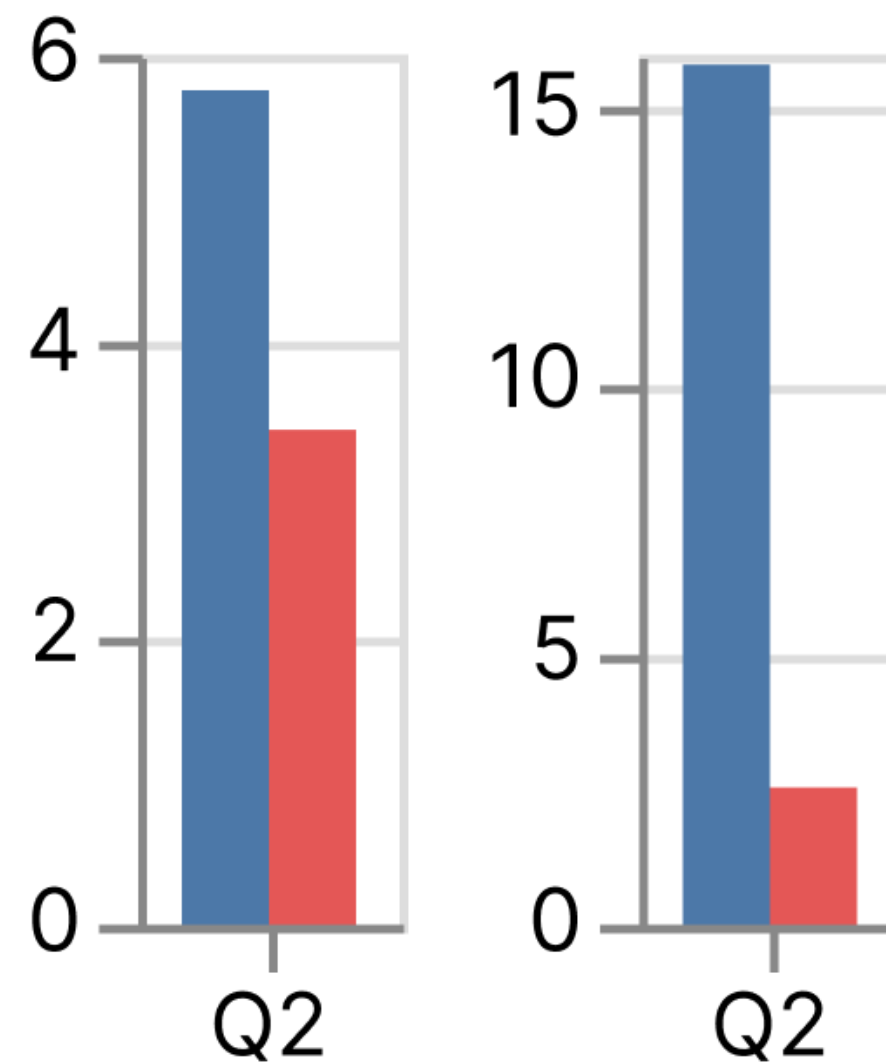
```
c = w.object(); p = w.object(); i = w.geoConstruct('intersection')
w.filter(c.type=='car' & p.type='people' & turnLeft(c)
        & contains(i, [c, p]))
```

Results



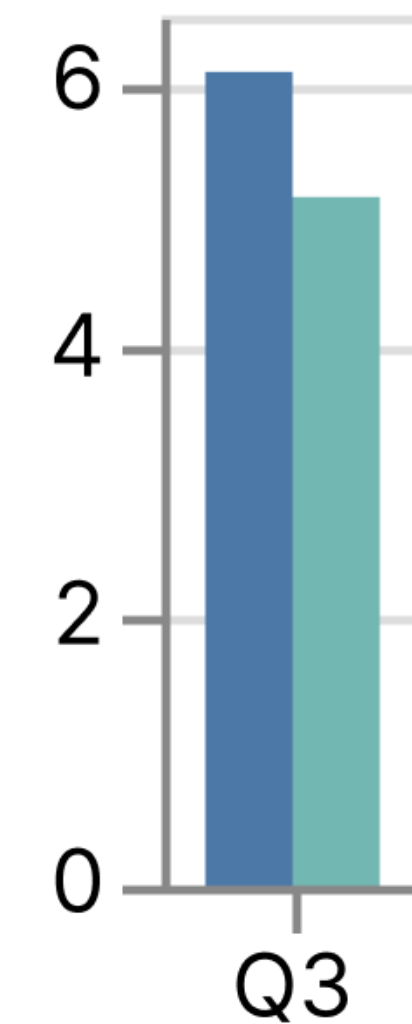
Ablation

Compare with no Optimization



VIVA

End-to-End video analytic system

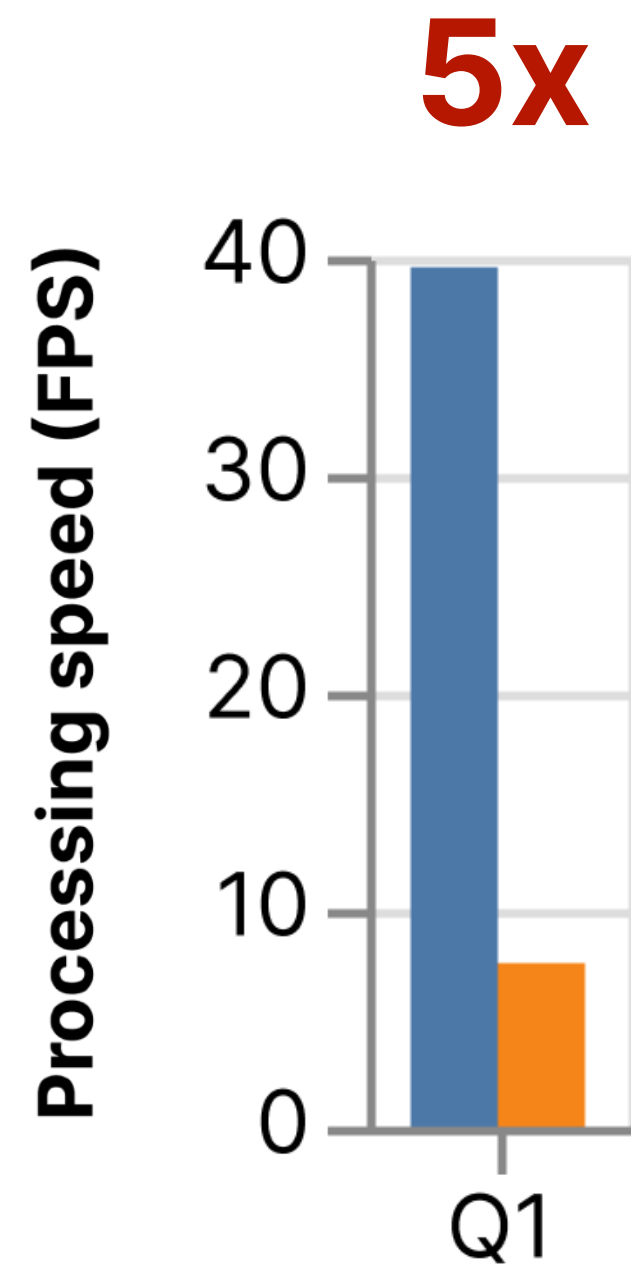


SkyQuery

Geospatial video query processor

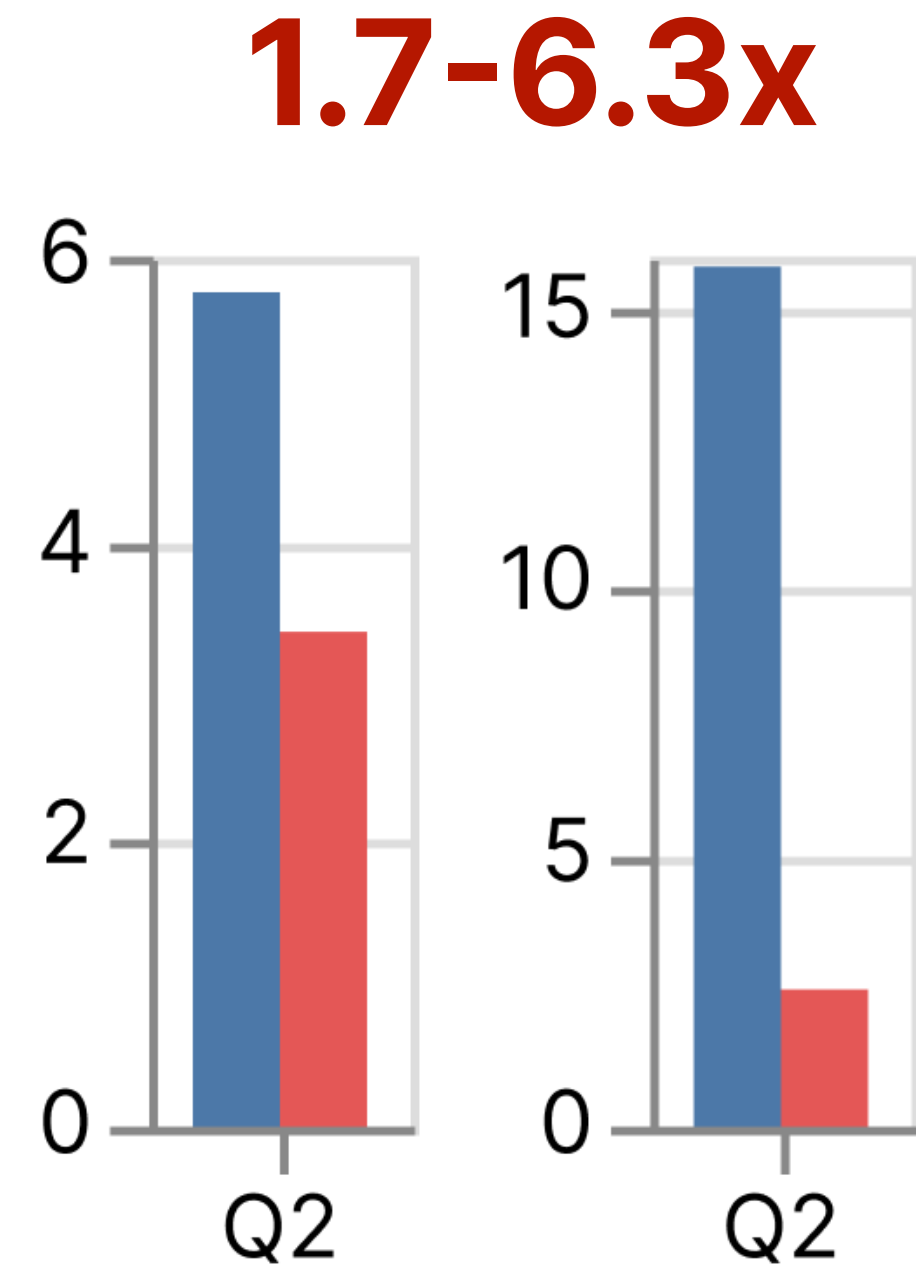
System: ■ Spatialyze ■ No Optimization ■ VIVA ■ SkyQuery

Results



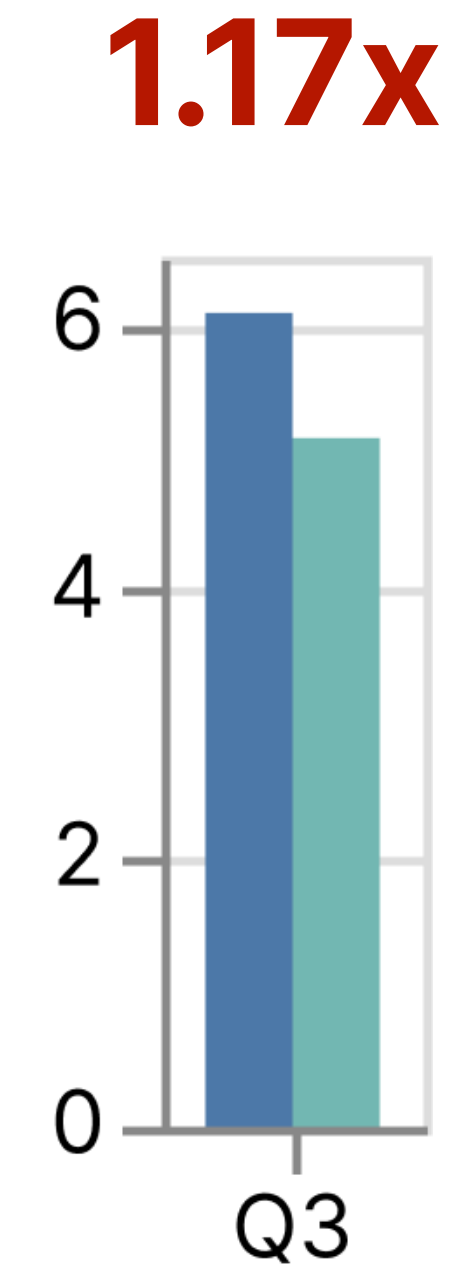
Ablation

Compare with no Optimization



VIVA

End-to-End video analytic system



SkyQuery

Geospatial video query processor

System: ■ Spatialyze ■ No Optimization ■ VIVA ■ SkyQuery

SPATIALYZE

In Summary



Constructing geospatial video analytics workflows is tedious and error-prone, let alone constructing them efficiently.



Spatialyze offers simplified data model and programming model; users analyze geospatial metadata together with videos.



Based on users' queries, Spatialyze accelerate video processing using geospatial metadata.

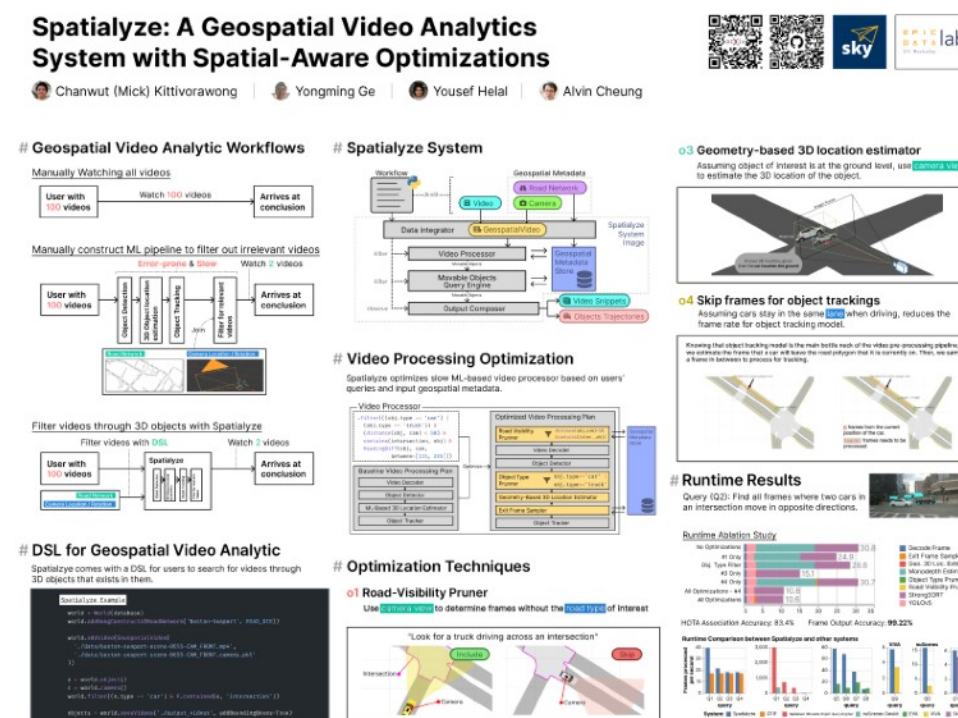
SPATIALYZE

Come visit us!



Code

github.com/apperception-db/spatialize



Poster

Tomorrow @ 10 AM



Paper

arxiv.org/abs/2308.03276

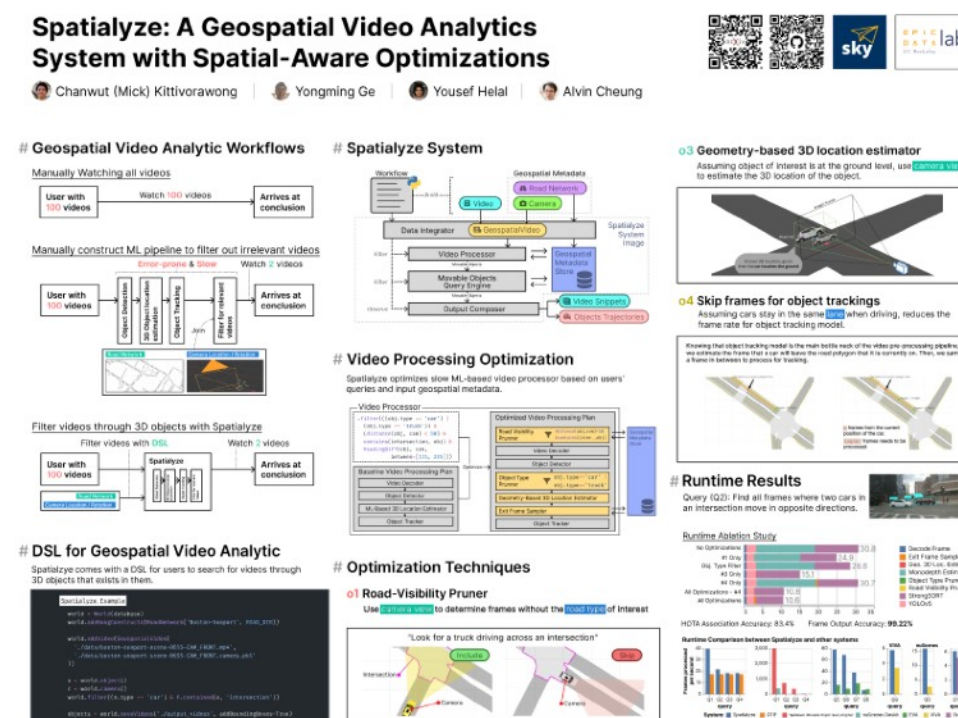
SPATIALYZE

Come visit us!



Code

github.com/apperception-db/spatialize



Poster

Tomorrow @ 10 AM



Paper

arxiv.org/abs/2308.03276

Next Step: User-Centered Design for Video Analytics Tools

- * We're starting a new research project!
- * Studying real-world video analysis workflows
- * Your insights will shape our design!

