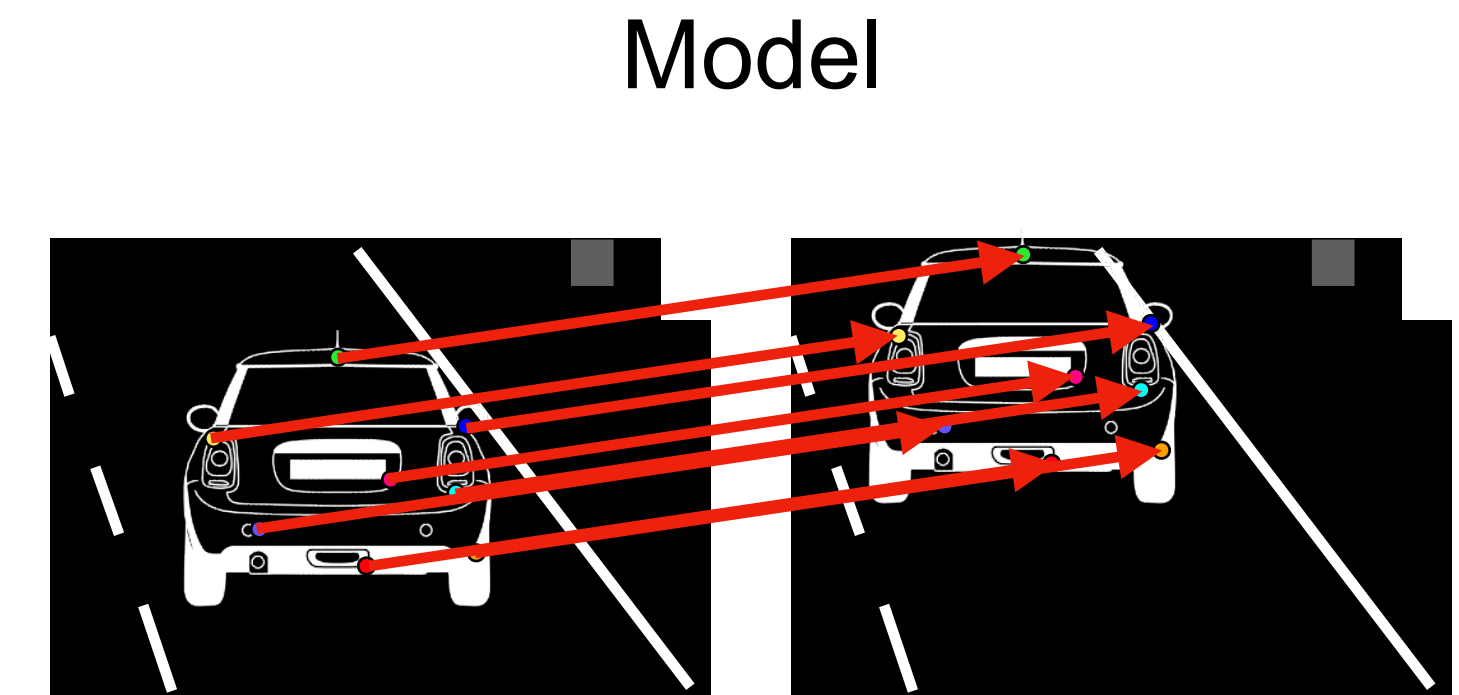
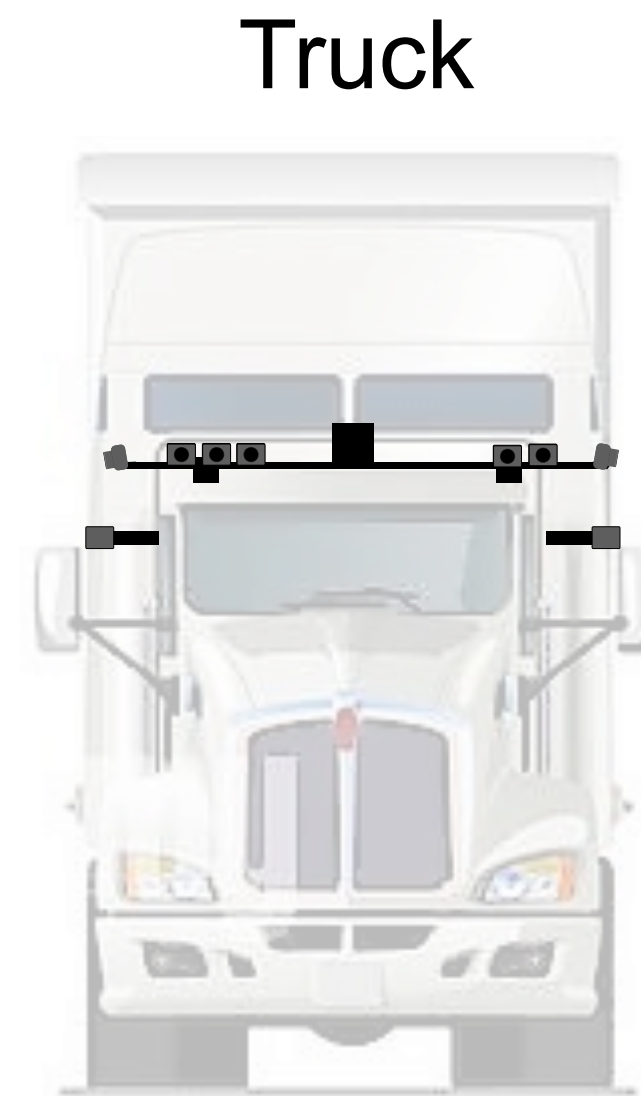


# Non-Player Character Visual Odometry

15 min

1. Introduction
2. Problem Constraints
3. Model
4. Implementation Details
5. Evaluation
6. Impact

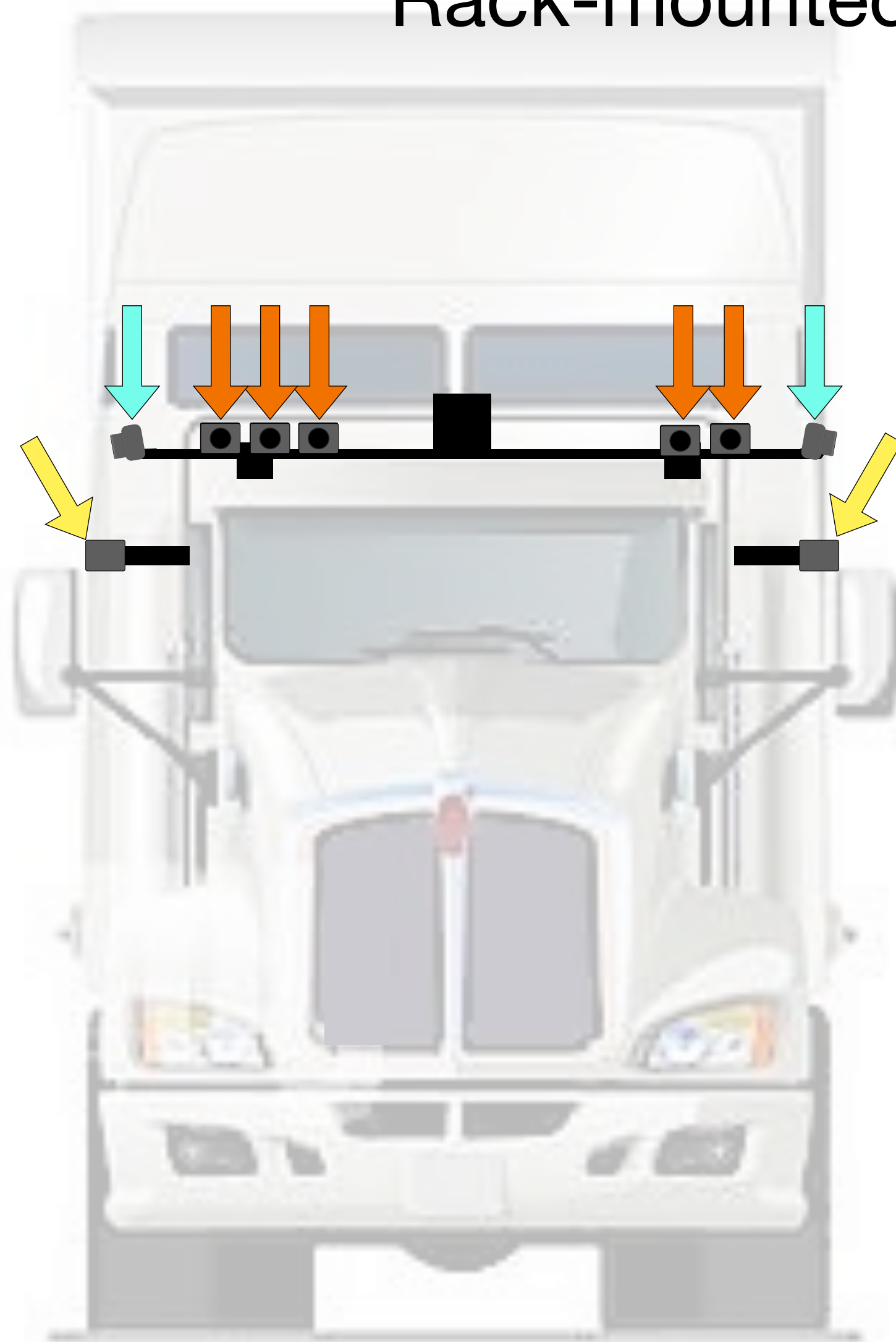


$$x^{t1} = K^{t1} H_{t0} K^{-1} x^{t0} + x_{corr}$$

Impact: Life-saving tech for autonomous vehicles.

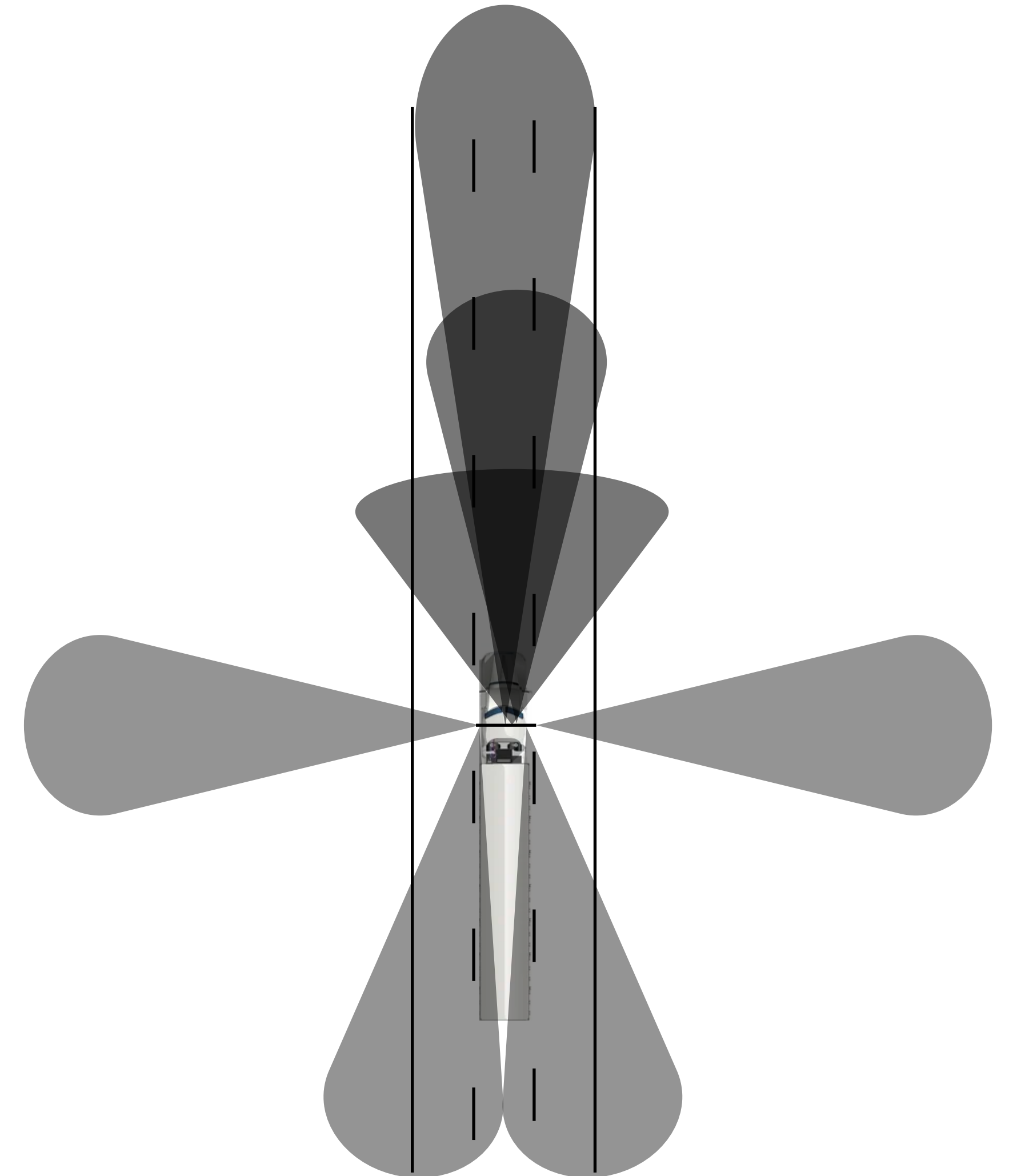
# 1. Introduction: 1. An Outfitted Truck

Rack-mounted cameras



- Front-facing
- Side-facing
- Rear-facing

Fields of view



# 1. Introduction: 2. Operational Design Domain

- L4 autonomous truck
- Highway, limited urban driving
- Mapped roads
- Routes are known a priori
- The same routes are repeatedly driven

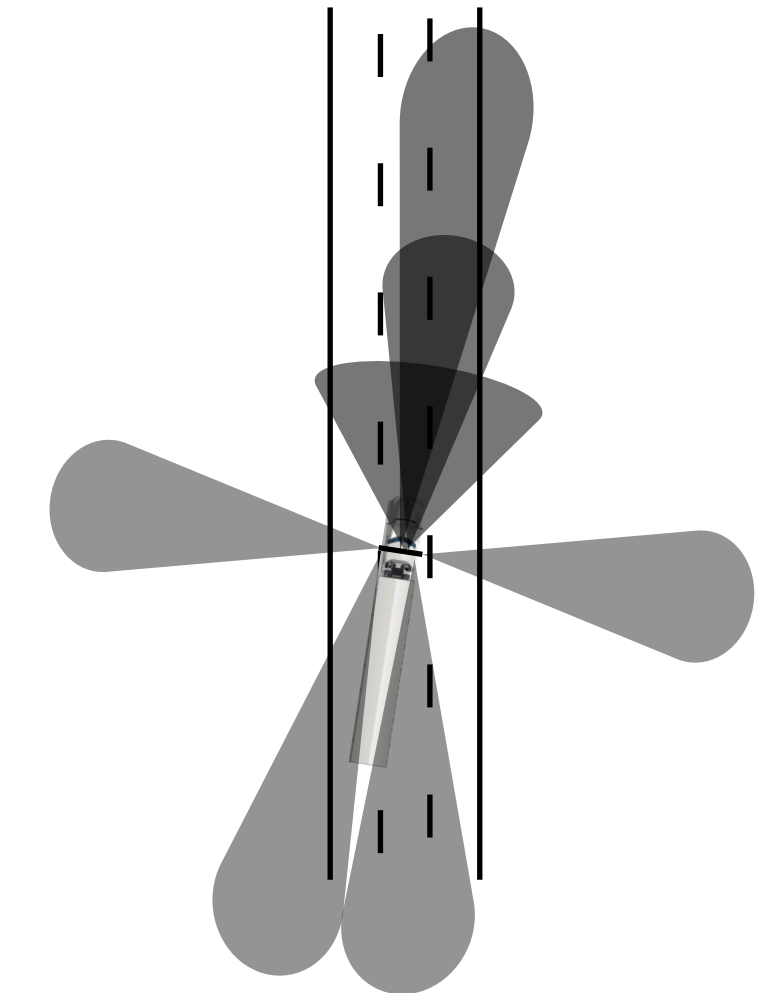
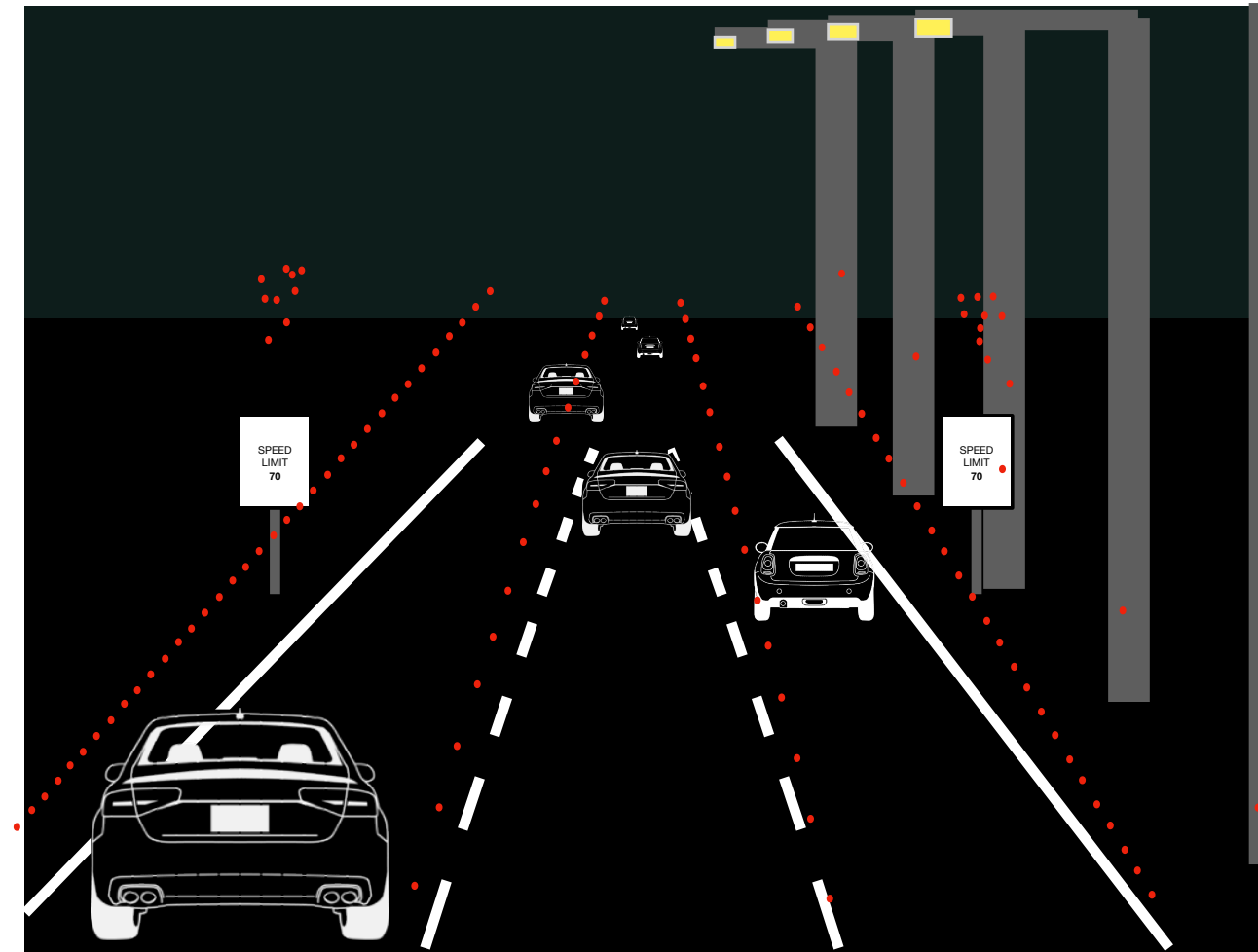


# 1. Introduction: 3. Camera Pose Estimation

Projection of an HD map onto an image

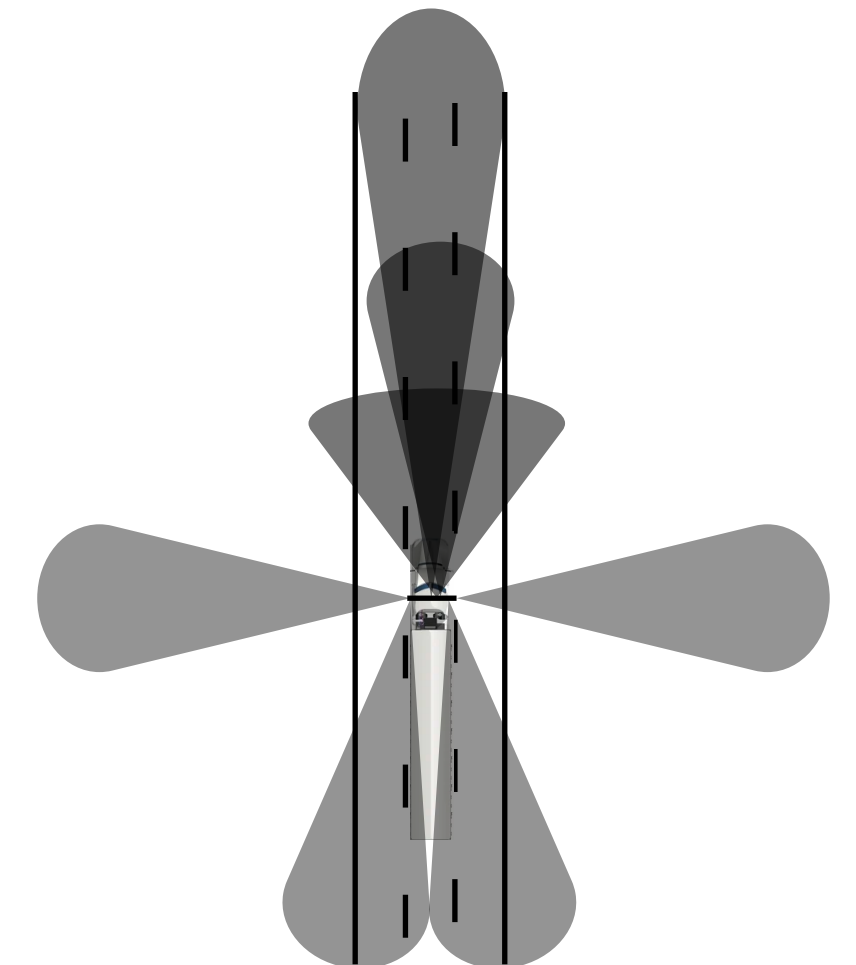
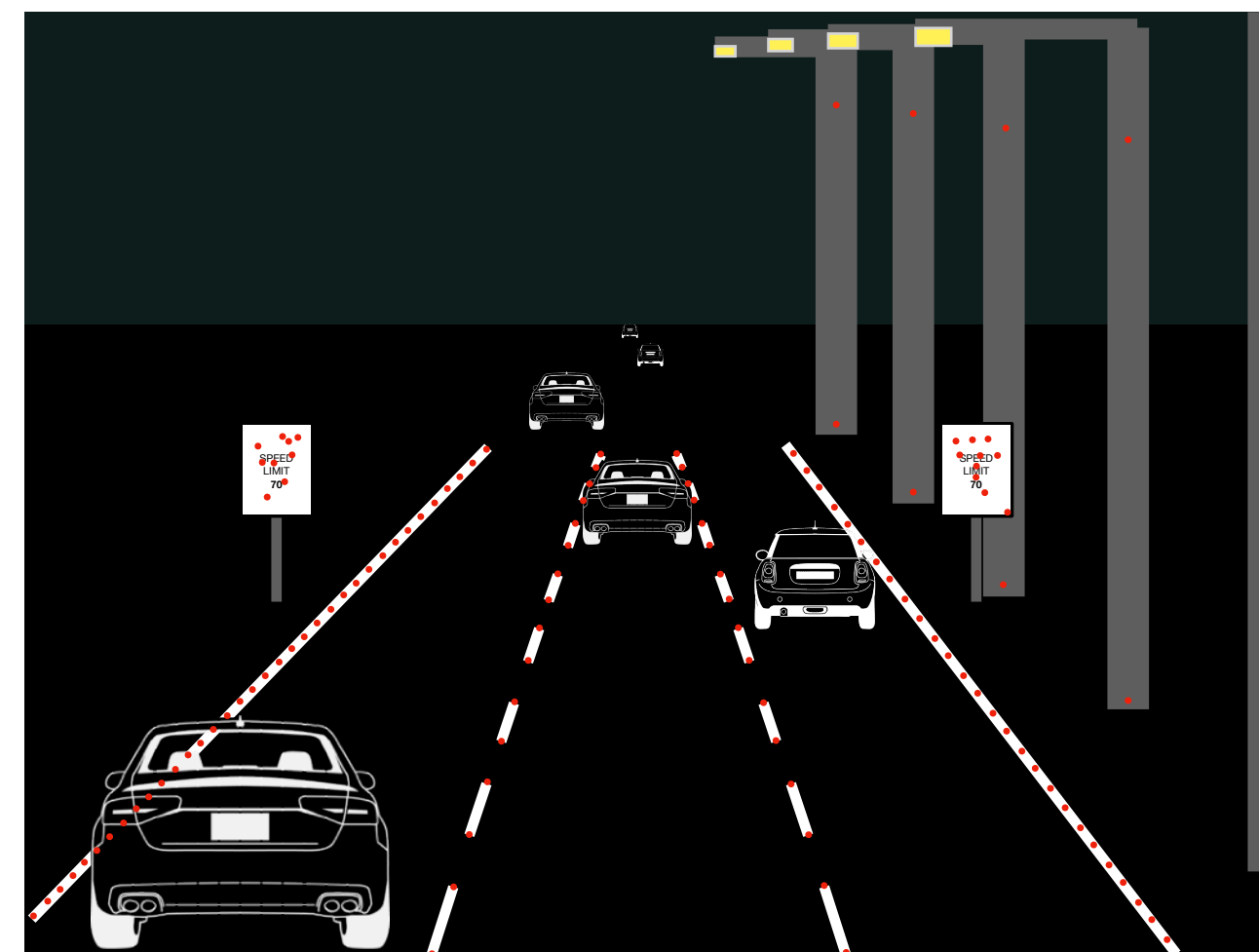
Camera FOVs

Unaligned  
(pitch and yaw error)



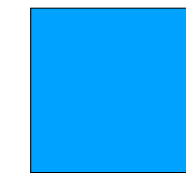
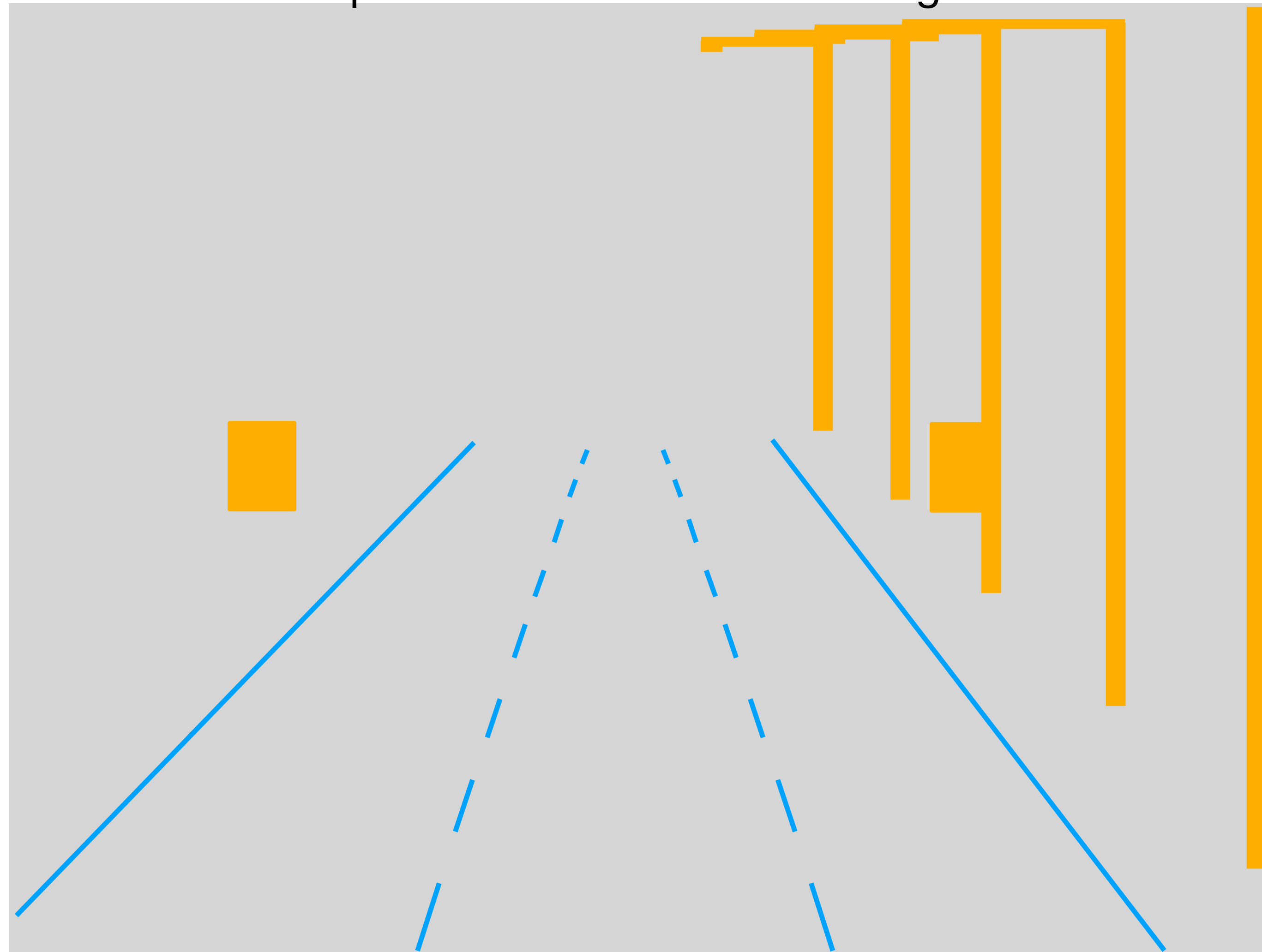
Goal: one  $SE(3)$  transform that brings all cameras into alignment

Aligned

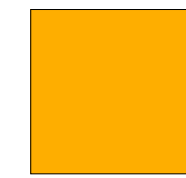


# 2. Problem Constraints: 1. Typical Visual Information

An example view from one front-facing camera



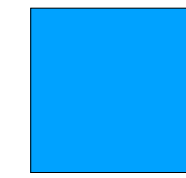
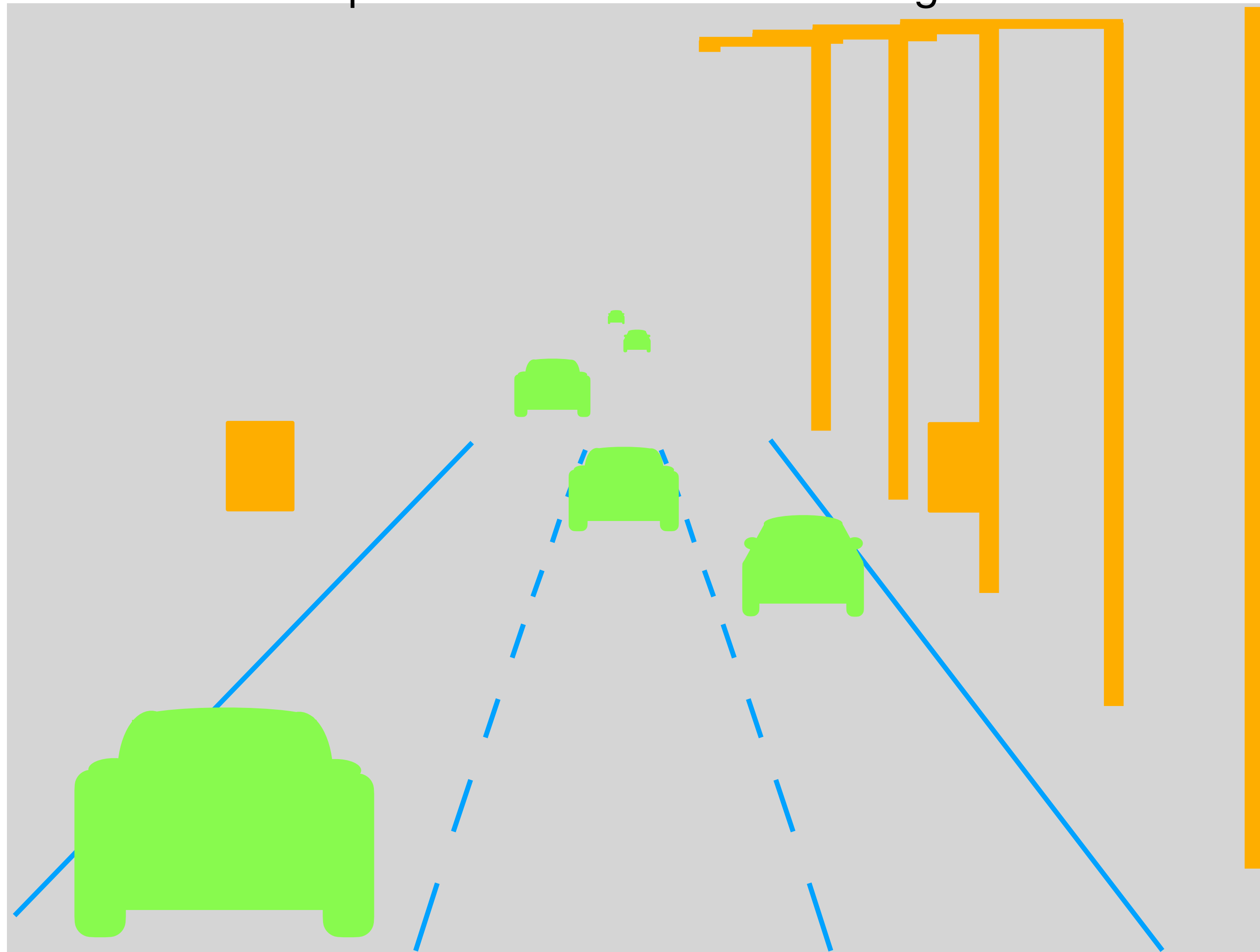
Primary



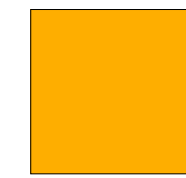
Secondary

# 2. Problem Constraints: 2. New Information Source

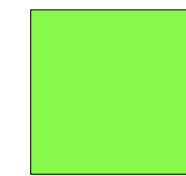
An example view from one front-facing camera



Primary



Secondary



This work

# 2. Problem Constraints: 3. Given Information

Camera

Non-Player Character (NPC)

Intrinsics\*

$$\mathbf{K} = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}$$

State

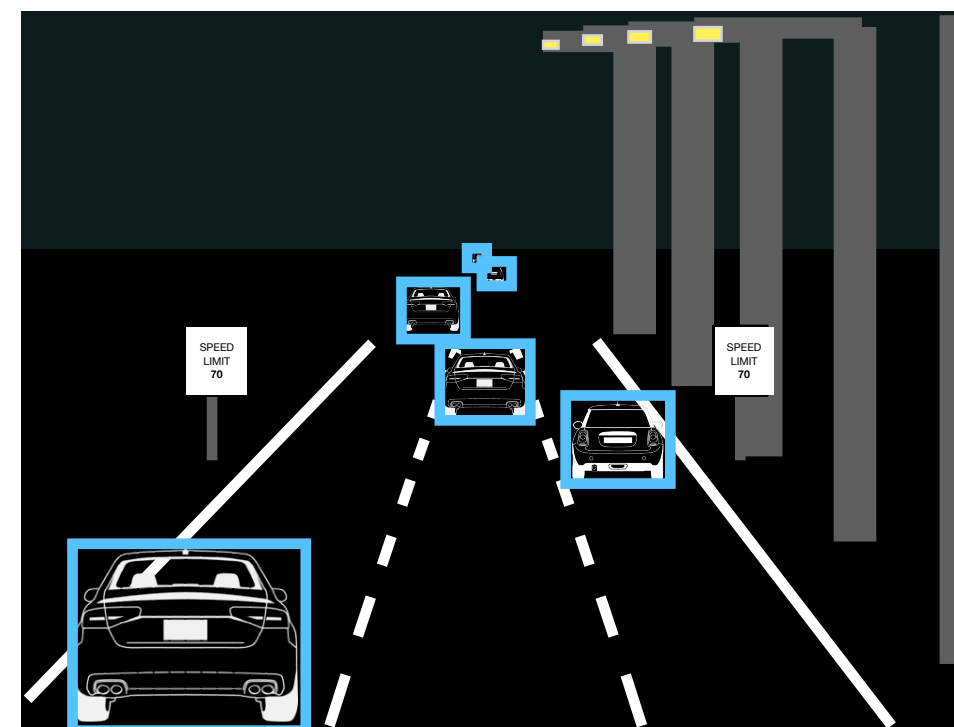
$$\mathbf{p}_{cam}^t = (x, y, z) \in \mathbb{R}^3$$

$$\mathbf{v}_{cam}^t = (\dot{x}, \dot{y}, \dot{z}) \in \mathbb{R}^3$$

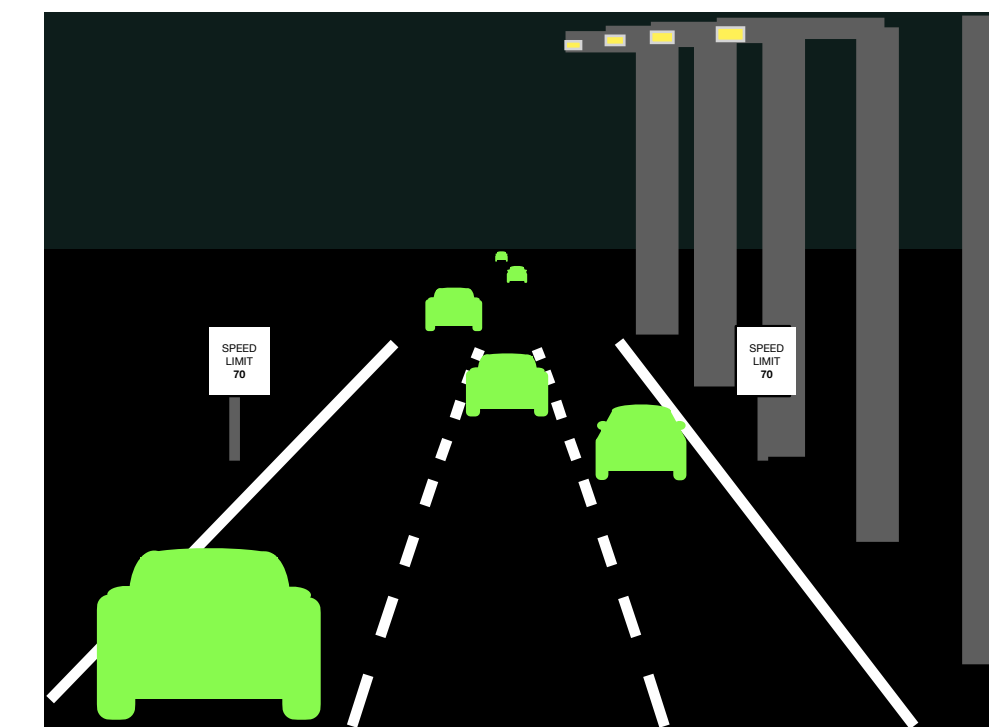
Identifier

128 byte REID

Bounding Box



Segmentation



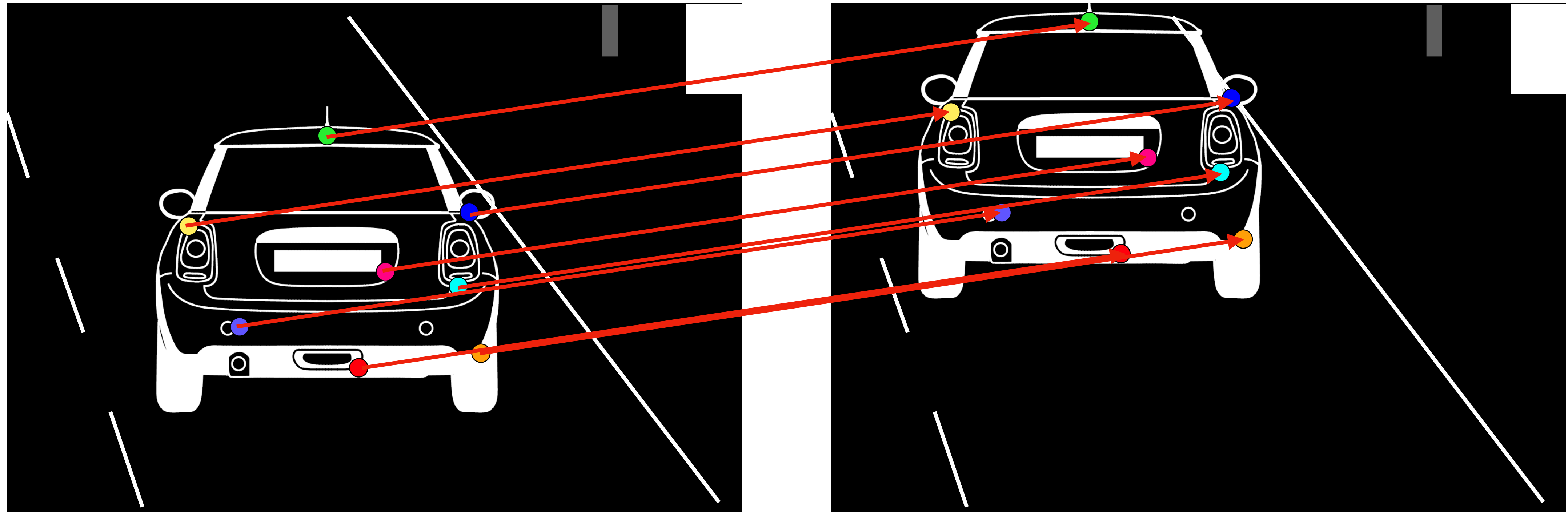
\*distortion is excluded to simplify this discussion



# 3. Model: 1. Where we're headed

Detected NPC keypoints,  $x^{t0}$

Projected NPC keypoints,  $\hat{x}^{t1}$ , from detections  $x^{t0}$



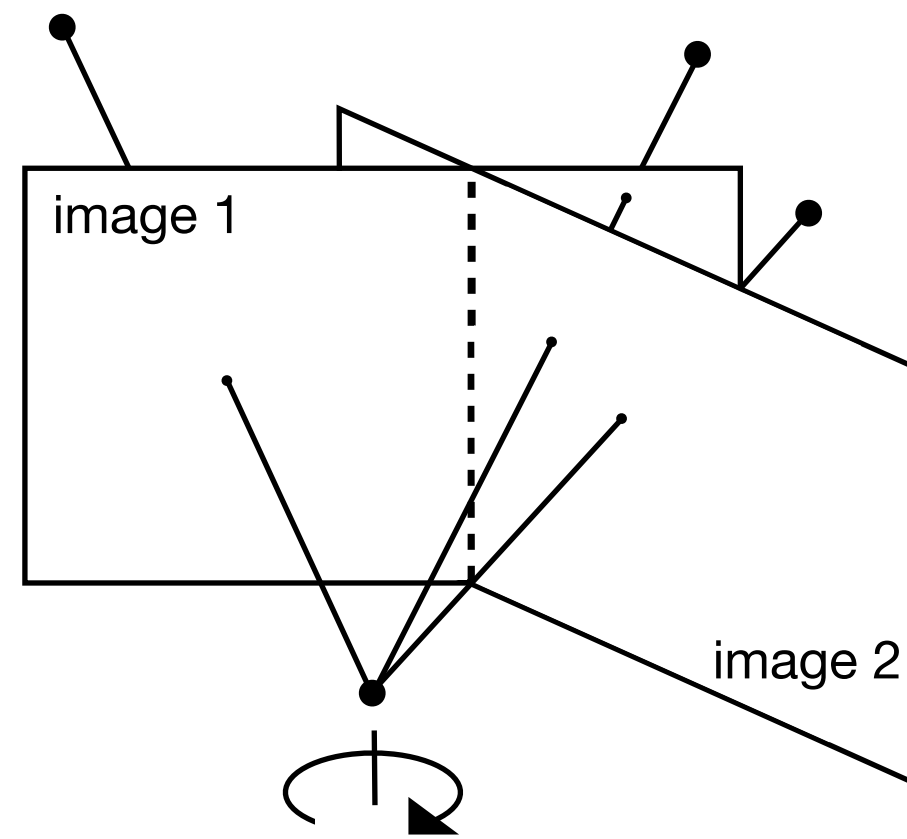
$$x^{t1} = K^{t1} H_{t0} K^{-1} x^{t0} + x_{corr}$$



# 3. Model: 2. Cases for the 2D Homography Transform

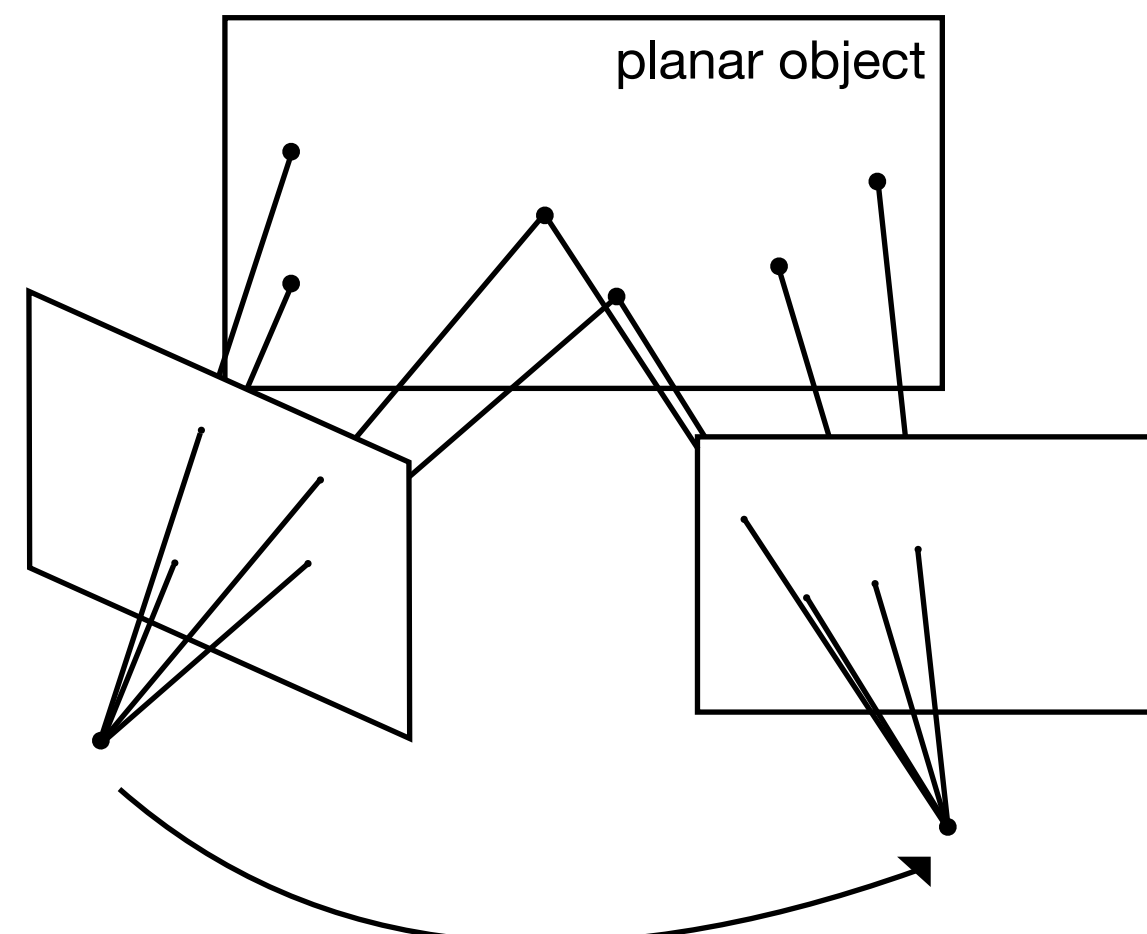
$$\hat{x}^{t1} = K {}^{t1}H_{t0} K^{-1} x^{t0}$$

Rotation-only camera motion



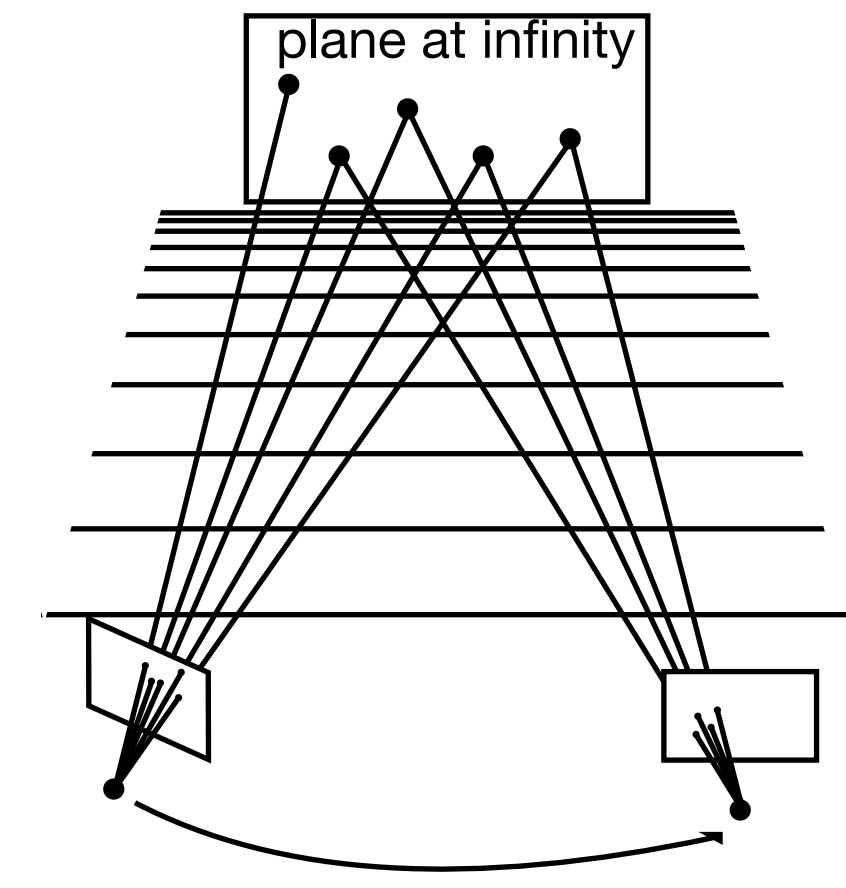
$${}^{t1}H_{t0} = {}^{t1}R_{t0}$$

Plane-induced homography



$${}^{t1}H_{t0} = {}^{t1}R_{t0} - \frac{n \cdot {}^{t1}t_{t0}}{d}$$

Points on the plane-at-infinity

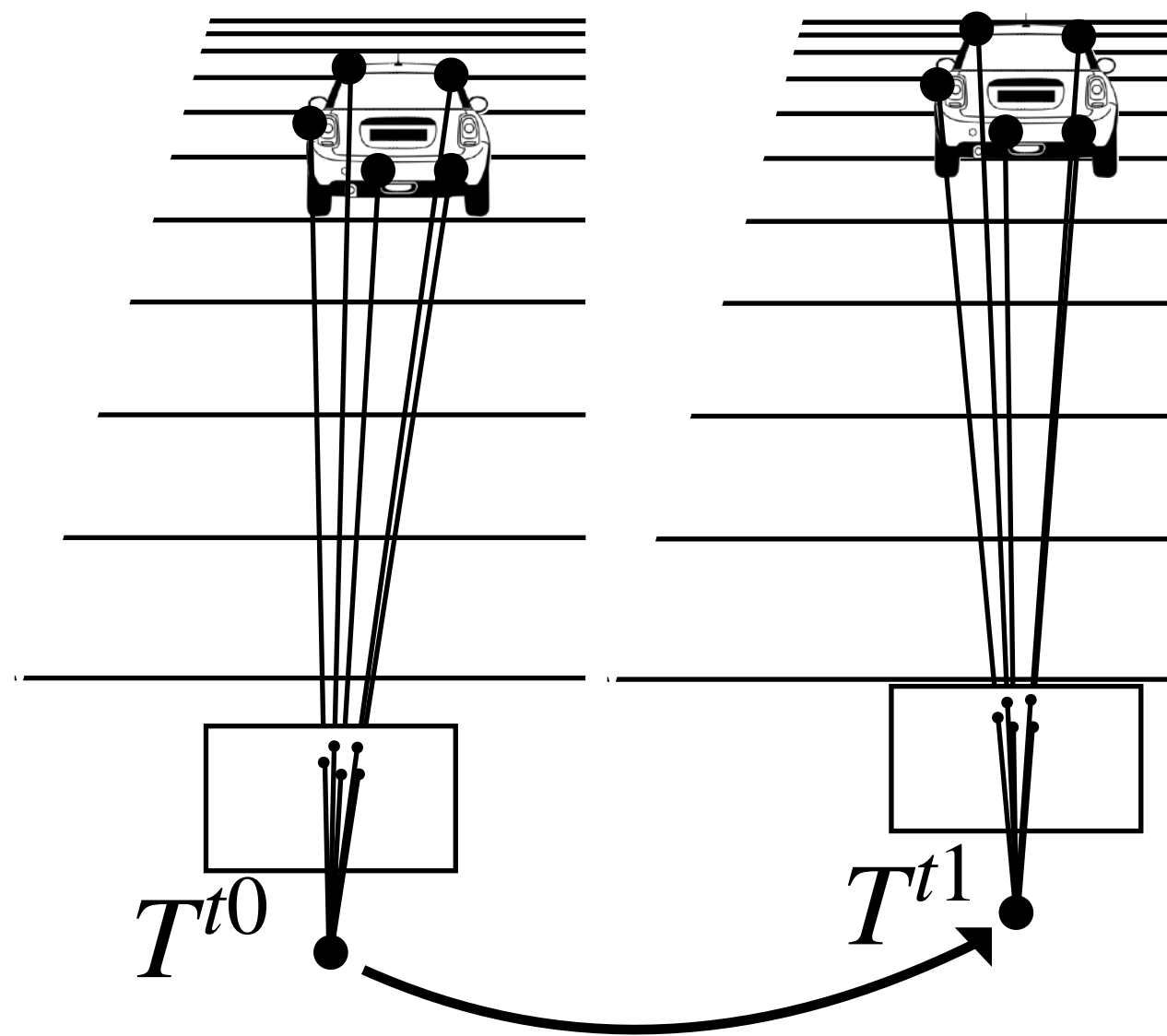


$${}^{t1}H_{t0} = {}^{t1}R_{t0}$$

# 3. Model: 3. Making the Plane-At-Infinity Assumption

$$\hat{x}^{t1} = K^{t1} H_{t0} K^{-1} x^{t0}$$

Plane-at-infinity assumption

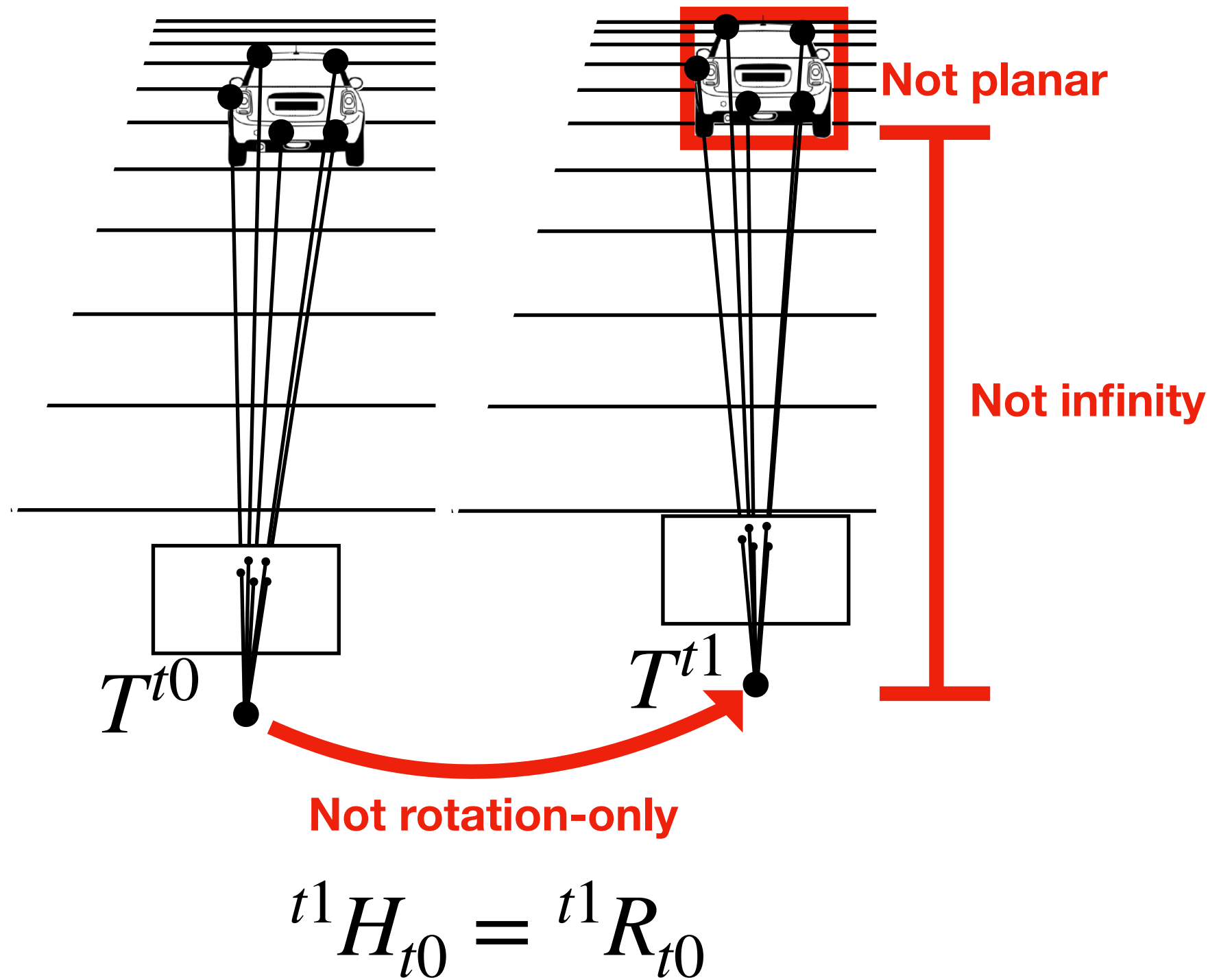


$${}^{t1}H_{t0} = {}^{t1}R_{t0}$$

# 3. Model: 3. Making the Plane-At-Infinity Assumption

$$\hat{x}^{t1} = K {}^{t1}H_{t0} K^{-1} x^{t0}$$

Plane-at-infinity assumption

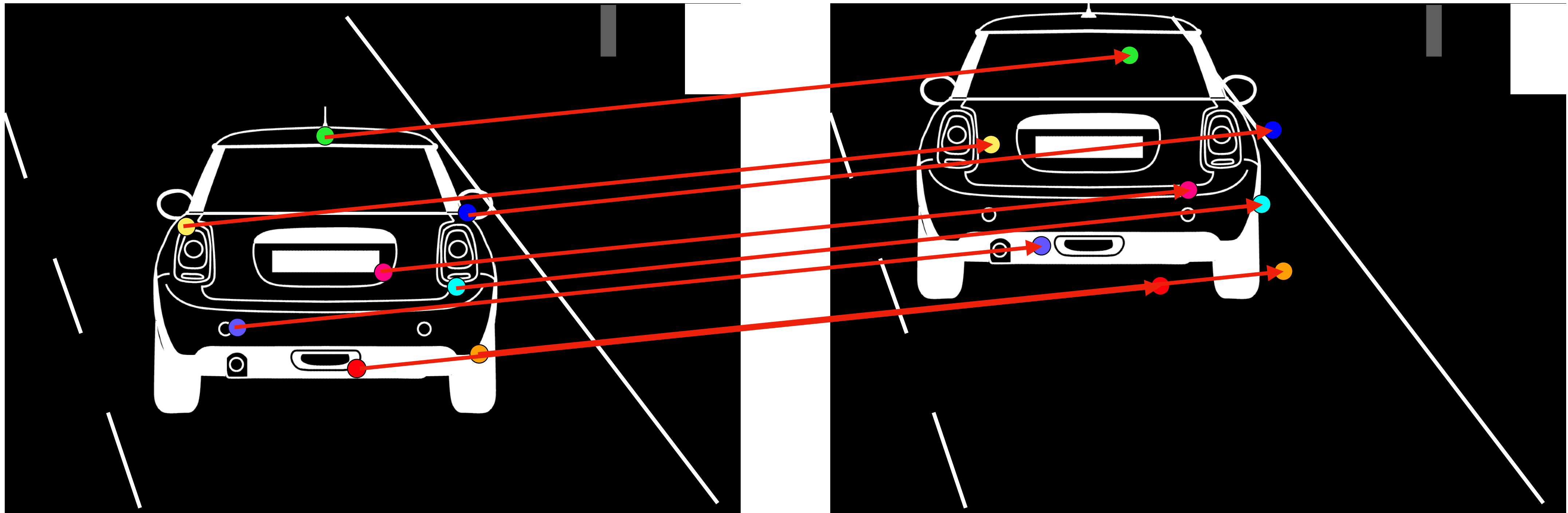


- But motion isn't rotation-only!
- NPCs are \*relatively\* stationary. They move with a similar velocity to the camera.
- Relative translation is small frame-to-frame.
- But NPCs are not at infinity!
- Many are close enough.
- But NPCs are non-planar!
- Non-planarity of NPCs becomes negligible farther away.

# 3. Model: 4. Kinematic Correction

Detected NPC keypoints,  $x^{t0}$

Projected NPC keypoints,  $\hat{x}^{t1}$ , from detections  $x^{t0}$

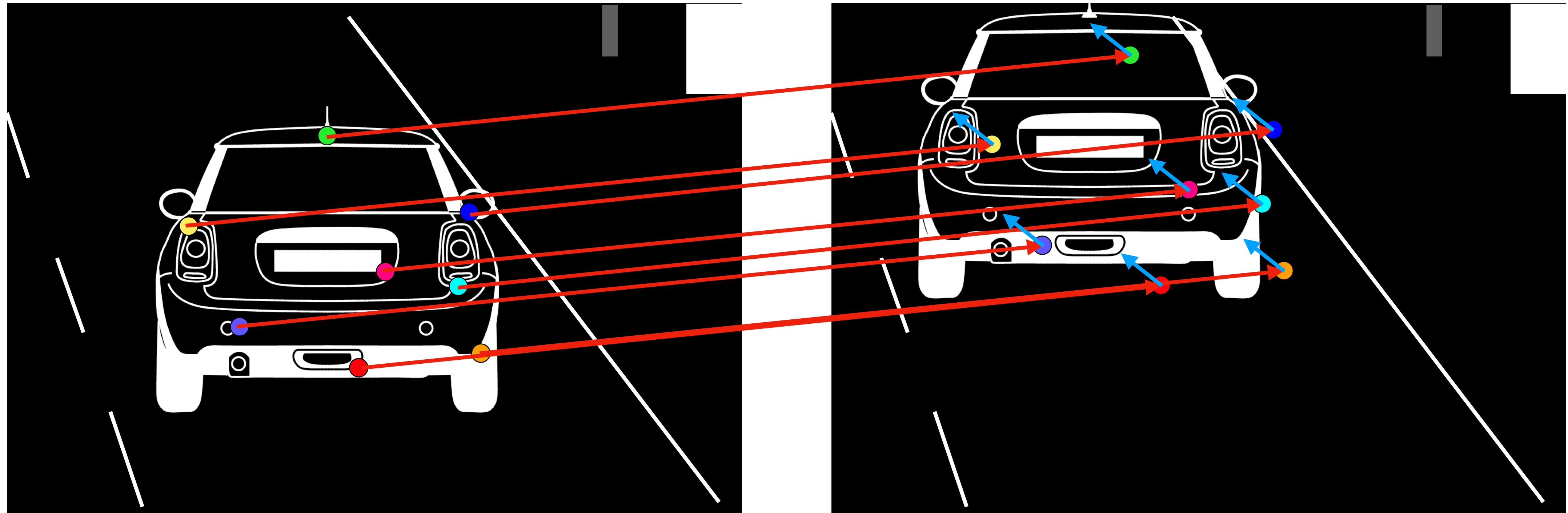


$$\hat{x}^{t1} = K^{t1} H_{t0} K^{-1} x^{t0}$$

# 3. Model: 4. Kinematic Correction

Detected NPC keypoints,  $x^{t0}$

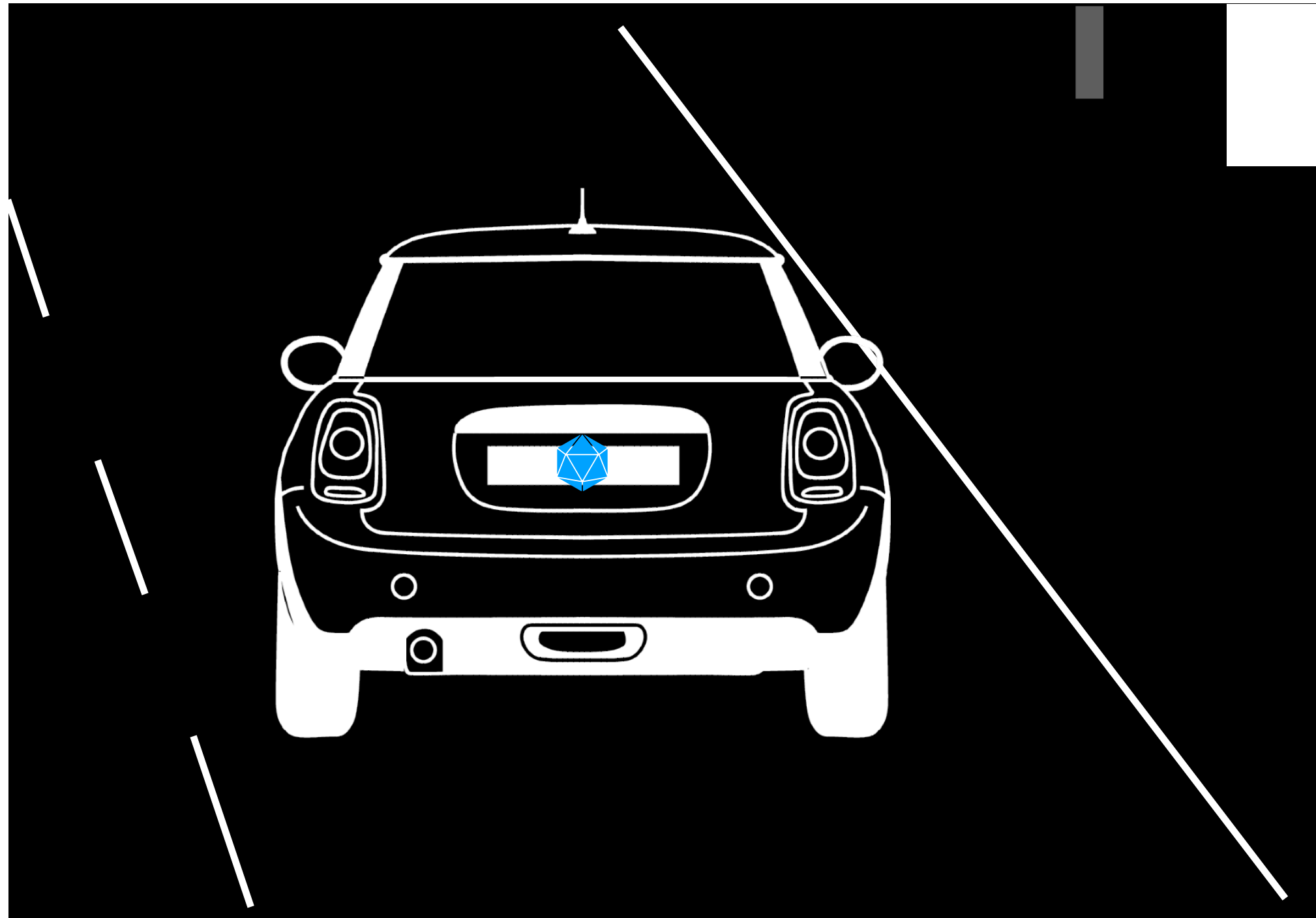
Projected NPC keypoints,  $\hat{x}^{t1}$ , from detections  $x^{t0}$



$$\hat{x}^{t1} = K^{t1} H_{t0} K^{-1} x^{t0} + \mathbf{X}_{\text{corr}}$$

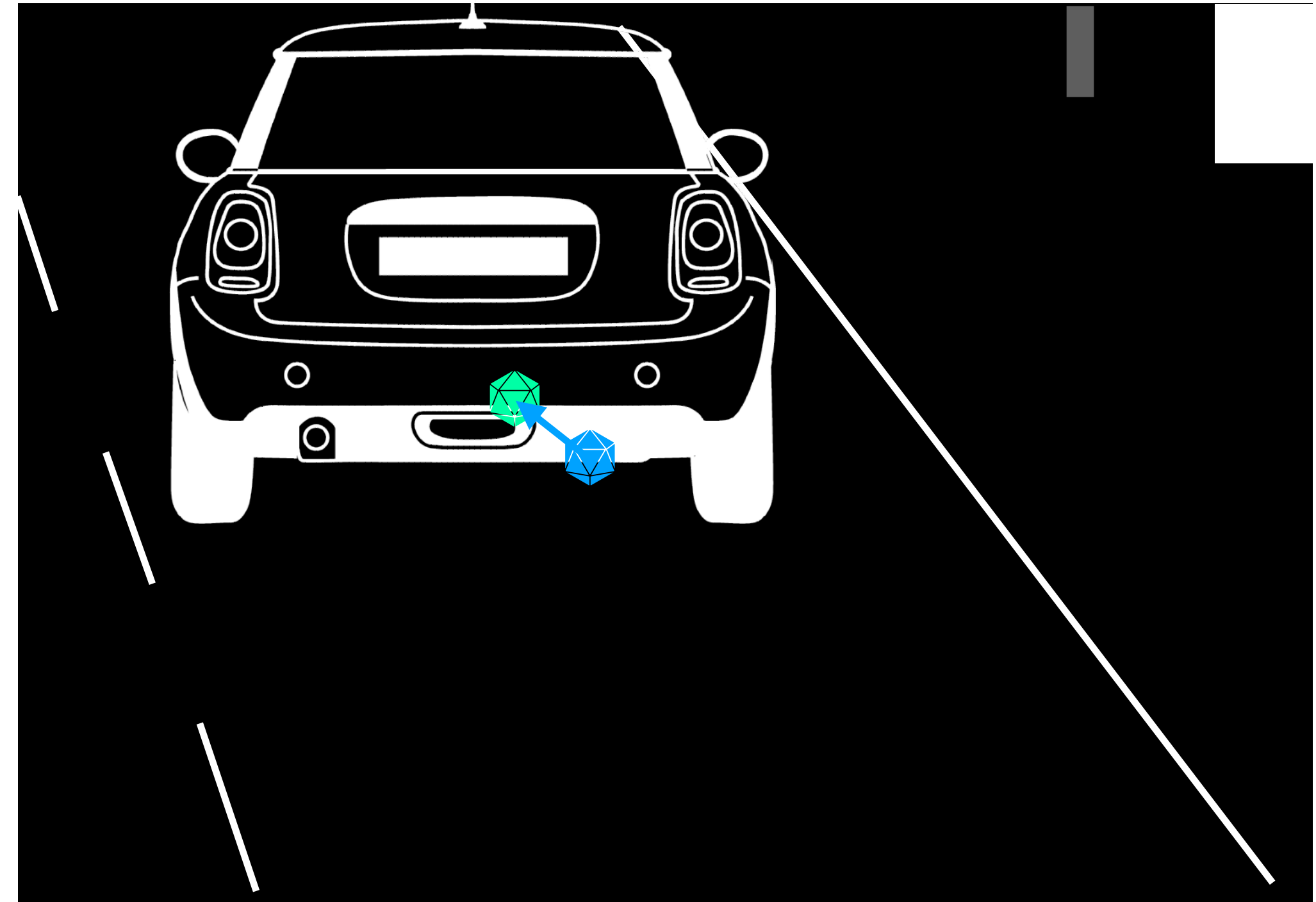
# 3. Model: 5. Obtaining the Kinematic Correction

Projected NPC position,  $p^{t_0}$



$$\text{blue dodecahedron } p^{t_0} = K p_{cam}^{t_0}$$

Projected NPC positions,  $p^{t_0}$  and  $\hat{p}^{t_1}$



$$\text{blue dodecahedron } p^{t_0}$$

$$\text{green dodecahedron } \hat{p}^{t_1} = K (p_{cam}^{t_0} + v_{cam}^{t_0} (t_1 - t_0))$$

$$\text{blue arrow } x_{corr} = \hat{p}^{t_1} - p^{t_0}$$



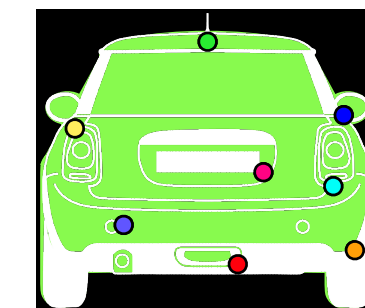
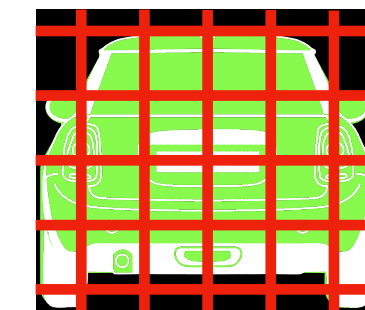
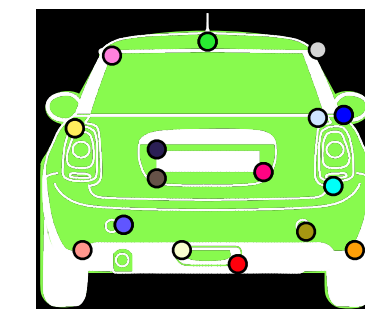
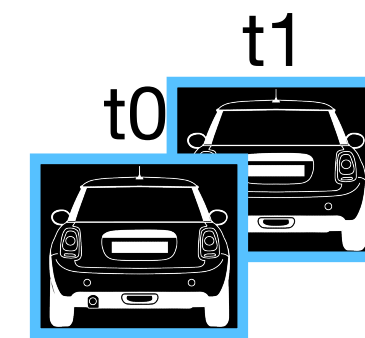
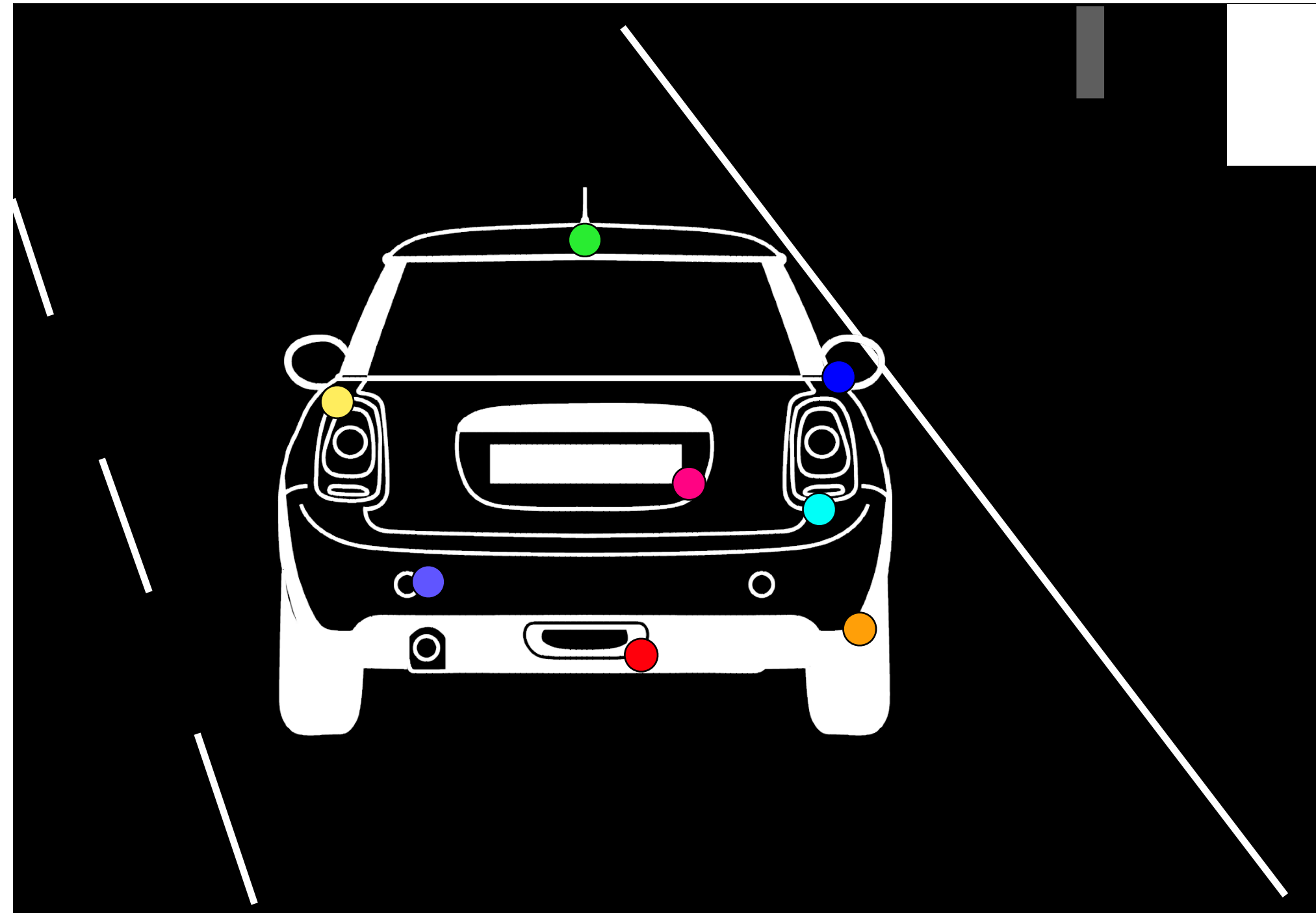
# 3. Model: 6. Use for Camera Pose Estimation

$$\operatorname{argmin}_{{}^{t1}\mathbf{R}_{t0}} \sum_{(i,j)} (x_j^{t1} - \pi(x_i^{t0}, {}^{t1}\mathbf{R}_{t0}))_{\Sigma}^2$$

${}^{t1}\mathbf{R}_{t0}$	$\text{SO}(3)$	orientation change between two camera frames
$(i, j)$	$\mathbb{N}^2$	indices of matched features in images at times t0 and t1
$x_i^{t0}$	$\mathbb{R}^2$	Pixel location of matched feature at time t0
$x_j^{t1}$	$\mathbb{R}^2$	Pixel location of matched feature at time t1
$\pi()$	$\mathbb{R}^2$	Our model equation: $\hat{x}^{t1} = K {}^{t1}H_{t0} K^{-1} x^{t0} + x_{corr}$
$()_{\Sigma}^2$	$\mathbb{R}^1$	We take the squared norm with respect to the covariance

# 4. Implementation Details: 1. Feature Extraction

Repeatable, matchable feature points



Steps

Matched NPCs from consecutive frames using REID

Extracted ORB feature points, within the segmentation boundary

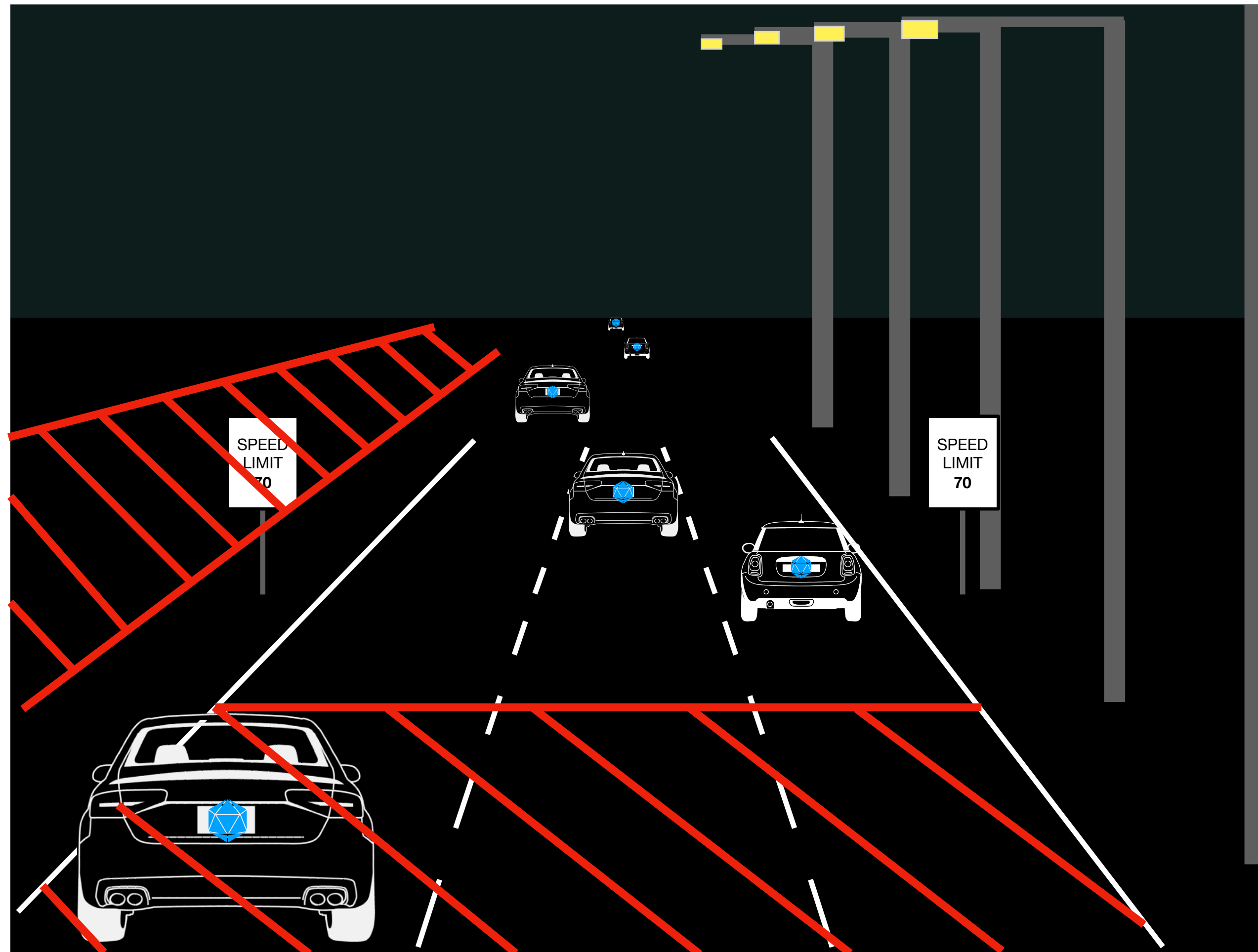
Improved coverage using a grid

Filter:

- match consistency
- homography
- accept sets with  $\geq 5$  points

\*\* Try ECC image alignment for higher fidelity

# 4. Implementation Details: 2. NPC Filtering



## Filters

- Closer than 75 m (this can vary per camera)
- Vehicles traveling in the opposite direction

75 m

# 5. Evaluation: 1. Main Results

- Residuals: Low alignment error.
- With respect to ground truth, which is obtained from sequences of well-aligned areas:
  - Pitch, yaw\*:  $0.2^\circ$
  - Roll:  $\sim 1-2^\circ$  (IIRC)
- NPC kinematic correction,  $x_{corr}$ , becomes unnecessary:
  - with more distance
  - on very straight roads
  - At speeds closer to relatively zero.

\*Possibility for very high fidelity visual odometry (better than  $0.2^\circ$ ) due to the sub-pixel data association accuracy on vehicles at 1000m+ range.

# 5. Evaluation: 2. Comparison of Information

	Lane markings	Traffic signs, poles	NPCs
<b>Information Type</b>	Localization	Localization	Visual Odometry
<b>Requires an HD Map</b>	Yes	Yes	No
<b>Map Overhead</b>	High	High	None
<b>Prone to mapping error</b>	Yes	Yes	No
<b>Requires NPC estimates</b>	No	No	Yes
<b>Visible Range</b>	Up to 250 m	1000m+	1000m+
<b>Occurrence</b>	Ubiquitous	Sporadic	Ubiquitous
<b>Image real-estate</b>	Mid. Sometimes occluded by NPCs	Typically low	Typically mid (standard lens) to prime (telephoto lens)
<b>Affected by perceptual aliasing</b>	Yes	No. Poles yes.	No
<b>Usable during sun glare</b>	If facing away from the sun.	Yes	Yes
<b>Night</b>	Visible if lit.	Visible if lit. Poles no.	Self-lit

# 6. Impact: 1. ADAS Modules

## Pose Estimation

- Extended ODD:
  - Works without a map: construction zones, unmapped roads and terrain.
  - Works despite bad road perception: rain, snow, golden hour sun-glare, Texas-faded lane marks.
- Gained Robustness: Extra information to eliminate outliers from other sources.
- Higher Fidelity: NPCs are visible beyond 1000 m, whereas lane mark are visible up to 250 m. Double the fidelity using data from both forward- and backward-facing cameras.
- Consistency: Pose estimates can be verified against the model.

## NPC Estimation

- Gained Robustness: Extra information of NPC kinematics in cases where  $x\_corr$  is non-zero.
- Consistency: NPC estimates can be verified against the model.



# 6. Impact: 2. ADAS Performance

- High-Fidelity Orientation Estimates: Crucial for accurate long-horizon path planning, e.g., 10+ second paths for merging and lane changes.
- Reduced Pitch Error: Prevents untimely hard-braking (distant NPC appears close) or unnecessary acceleration (close NPC appears distant).
- Improved Lane-Keeping: Accurate heading ensures corrections keep the vehicle on the planned path.
- Enhanced NPC Response: More accurate NPC kinematics improve adaptive responses to NPC behavior.
- Increased Passenger Comfort: Smoother transitions in speed and steering adjustments.
- Accelerated ADAS Deployment: Faster time-to-market for life-saving autonomous technology.

# 6. Impact: 3. Other Domains?

- Air-to-air?
- Surface-to-surface?
- Pedestrians?
- How General?

# Questions?