

Multi-view Pedestrian Detection

Aung Sithu

MS-Student

12.10.23

Contents

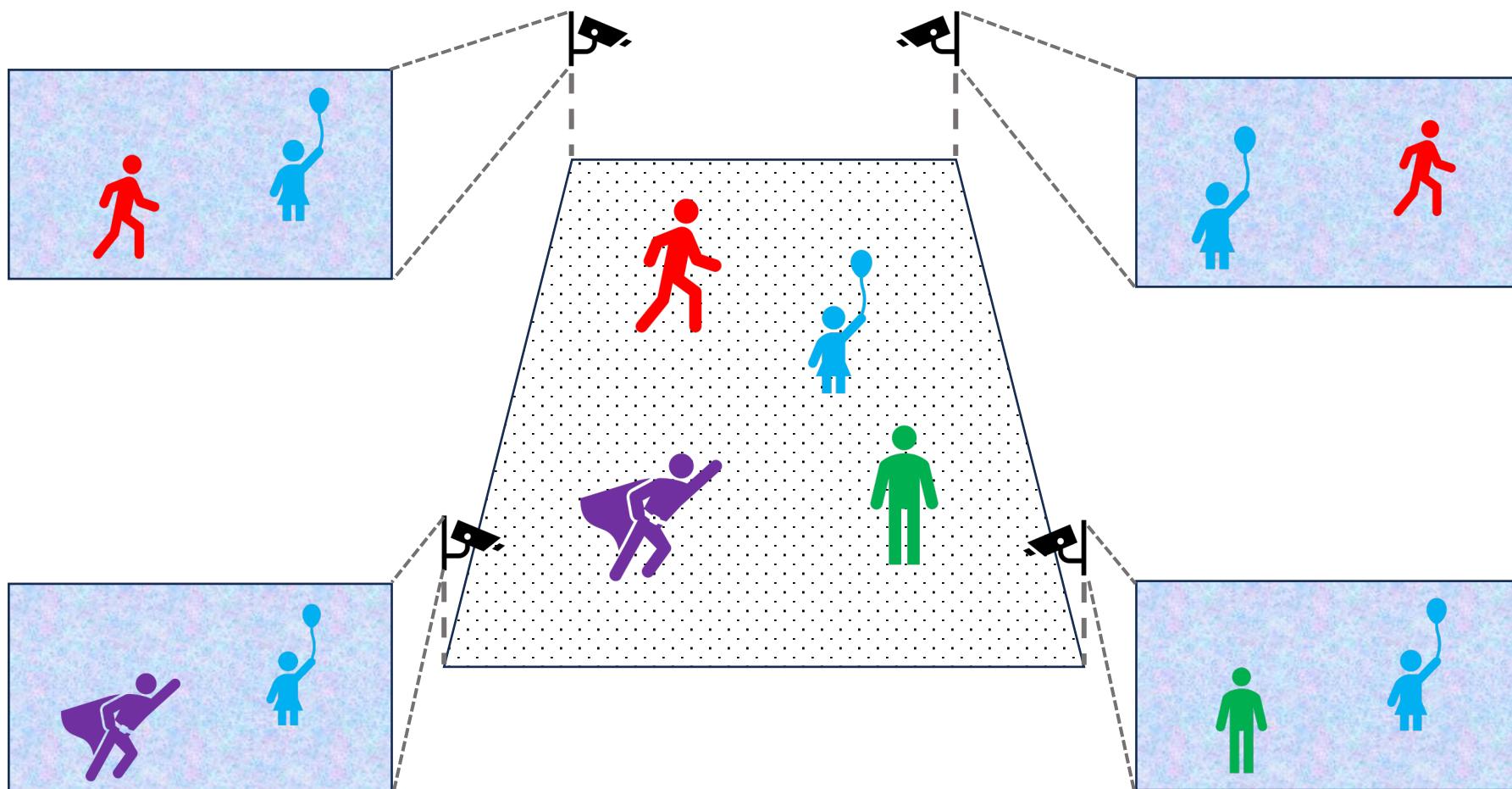
- Background
- Preliminaries
- Related Work (2D to BEV Feature Transform)
- Previous Work (Multi-view Detection)

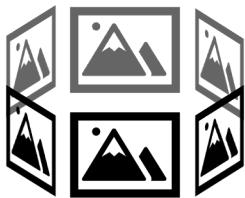


Background



Multi-camera Indoor People Localization

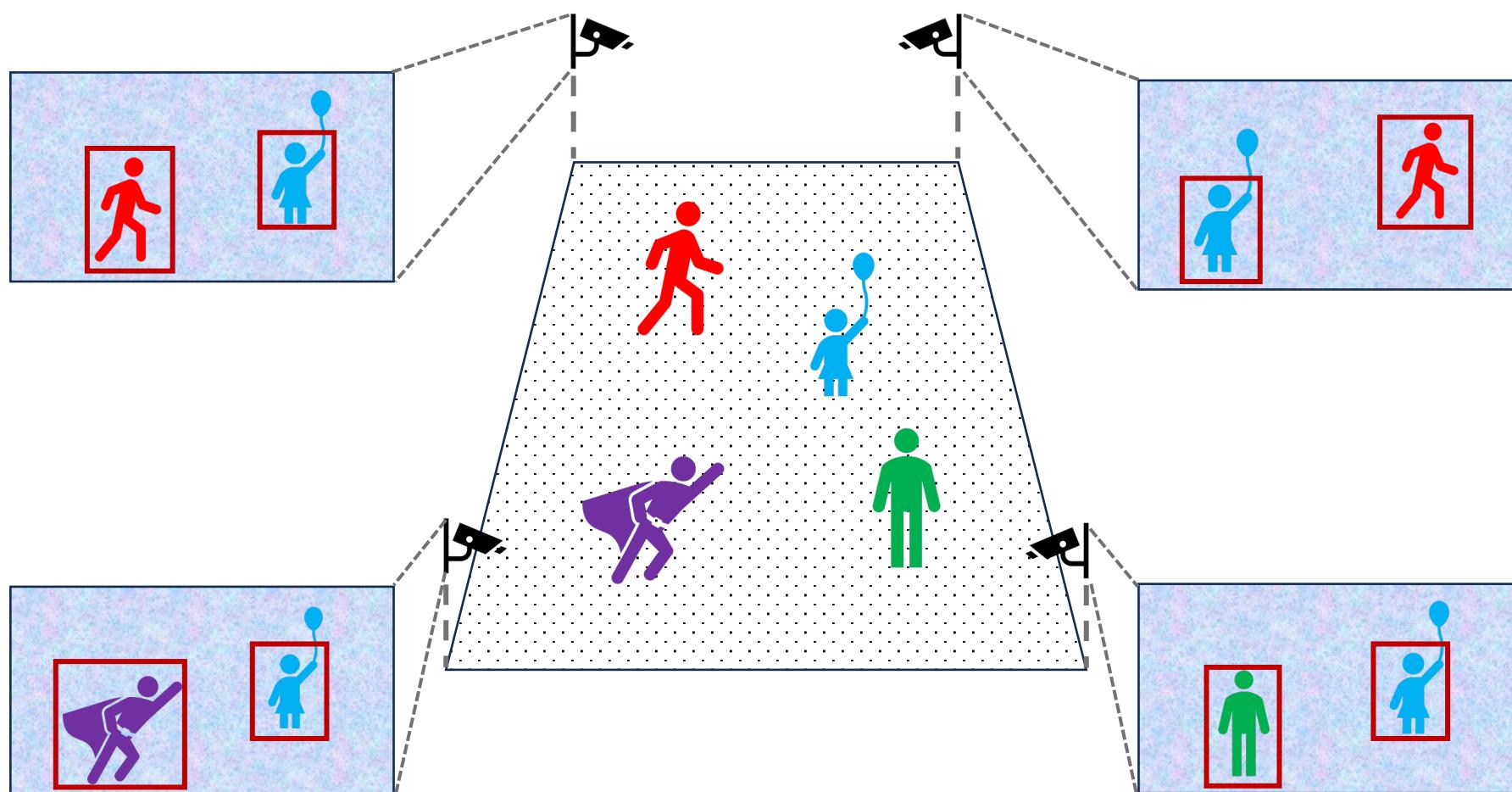


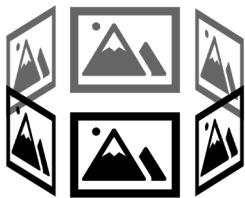


Background



Conventional Approach (2D Detection + Similarity Ranking + MOT)

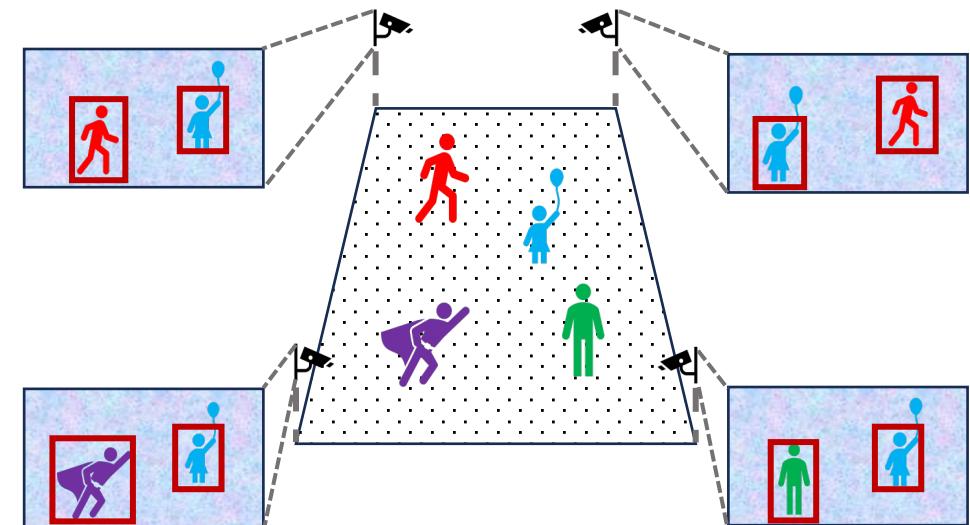
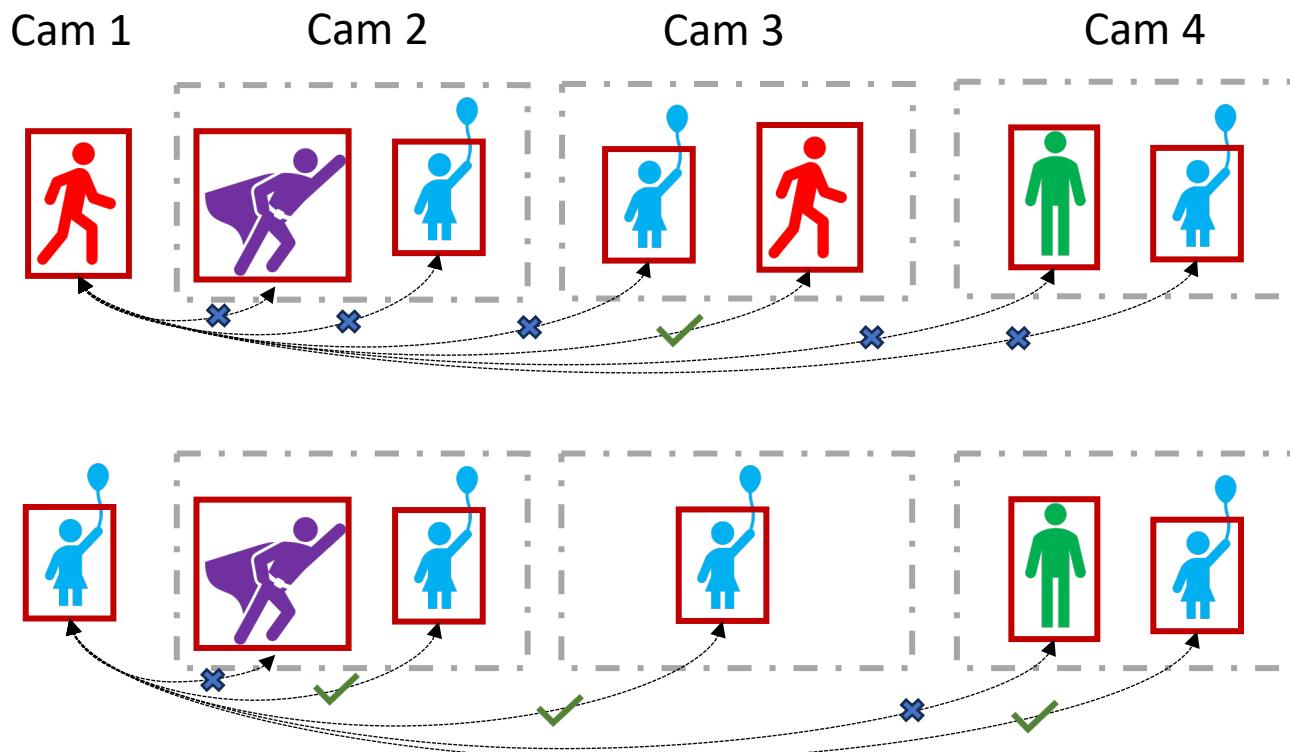




Background



Conventional Approach (Continued)



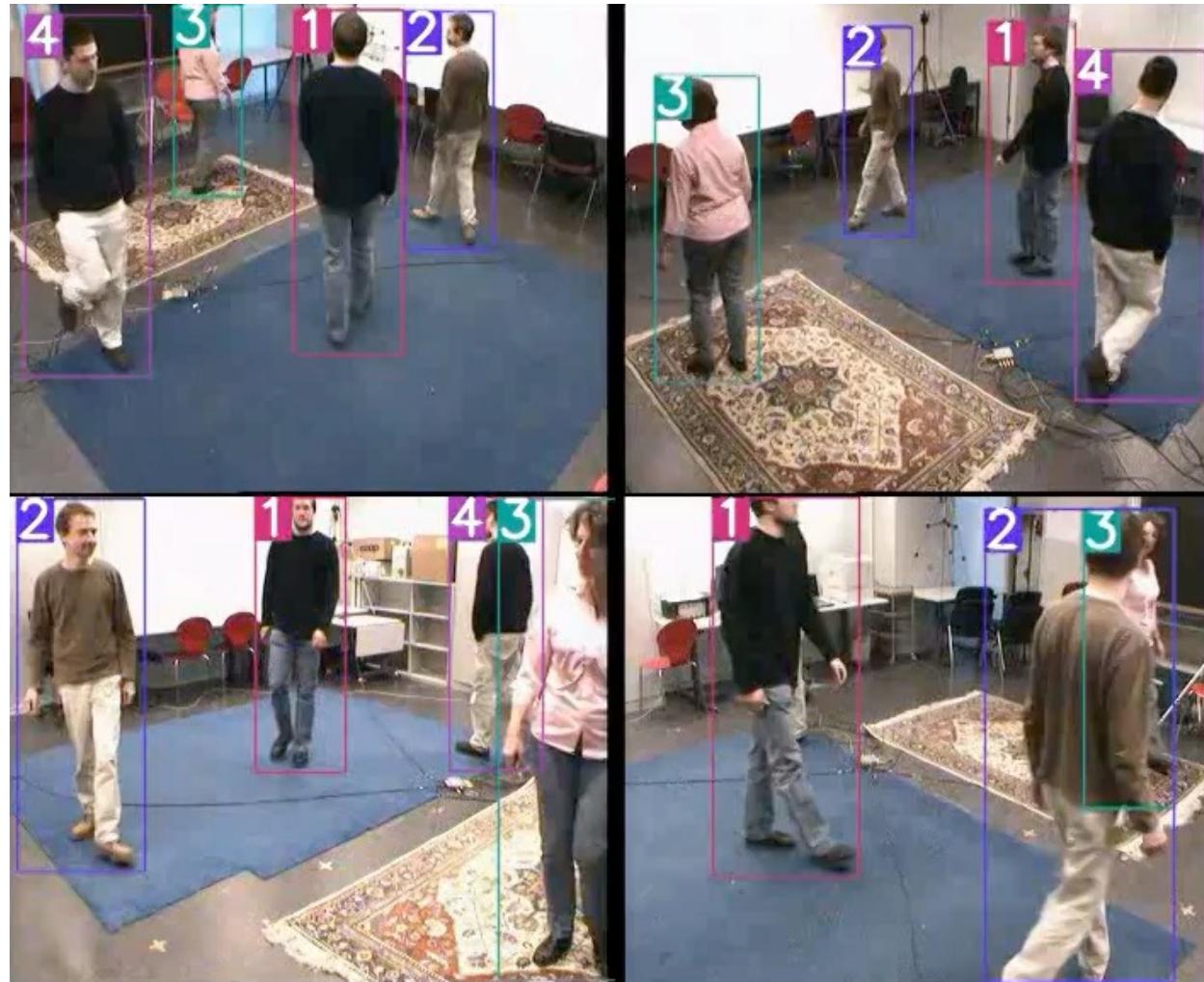
Severe Occlusions



Background



Conventional Approach (Continued)



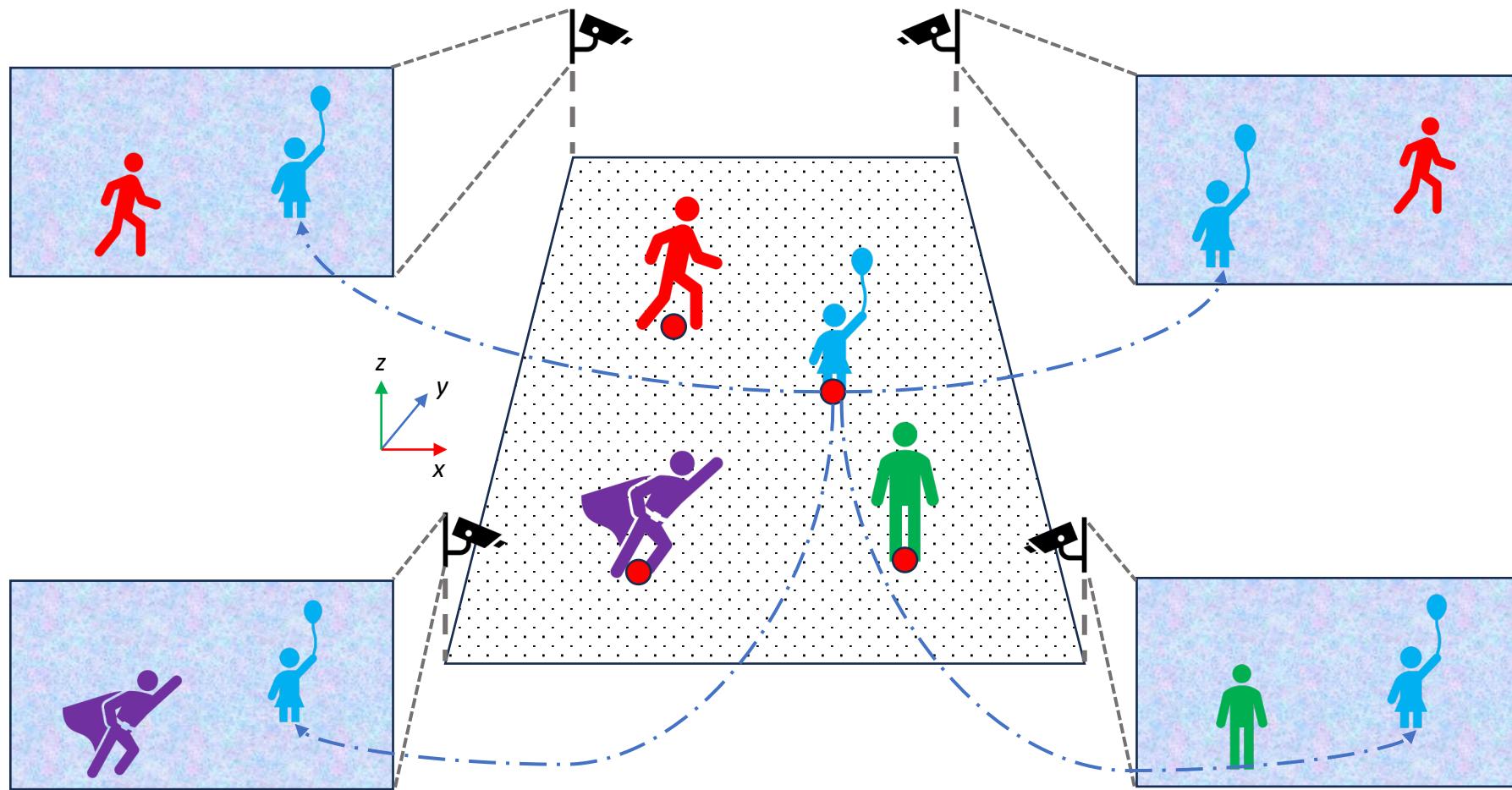


Background



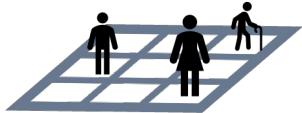
Towards 3D Bird's-Eye-View (BEV) Perception

Calibrated Cameras

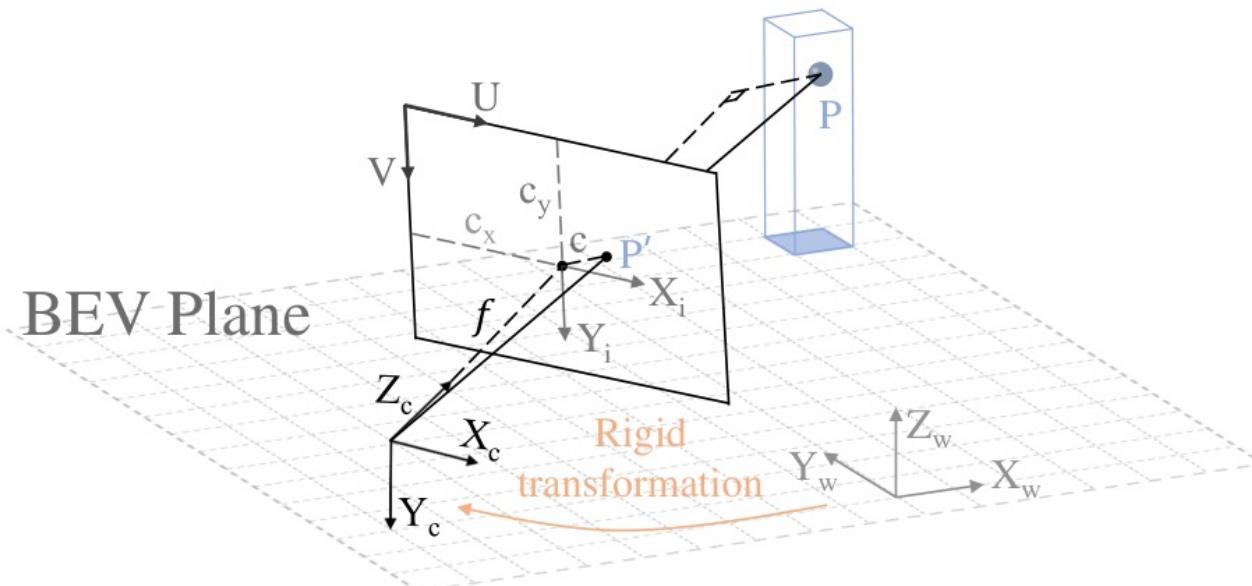




Preliminaries



Projective Geometry



Given the
world
coordinate of P

$$K[R \ T]$$

Retrieve the pixel
coordinates of P'

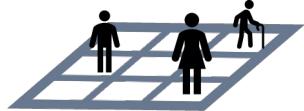
- $P = (X_w, Y_w, Z_w)$ -> world coordinate
- $P' = (X_c, Y_c, Z_c)$ -> camera coordinate
- $P' = (X_i, Y_i)$ or (U, V) -> pixel coordinate or image coordinate

$$z \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R \ T] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

- K = intrinsic
- $[R \ T]$ = extrinsic



Preliminaries



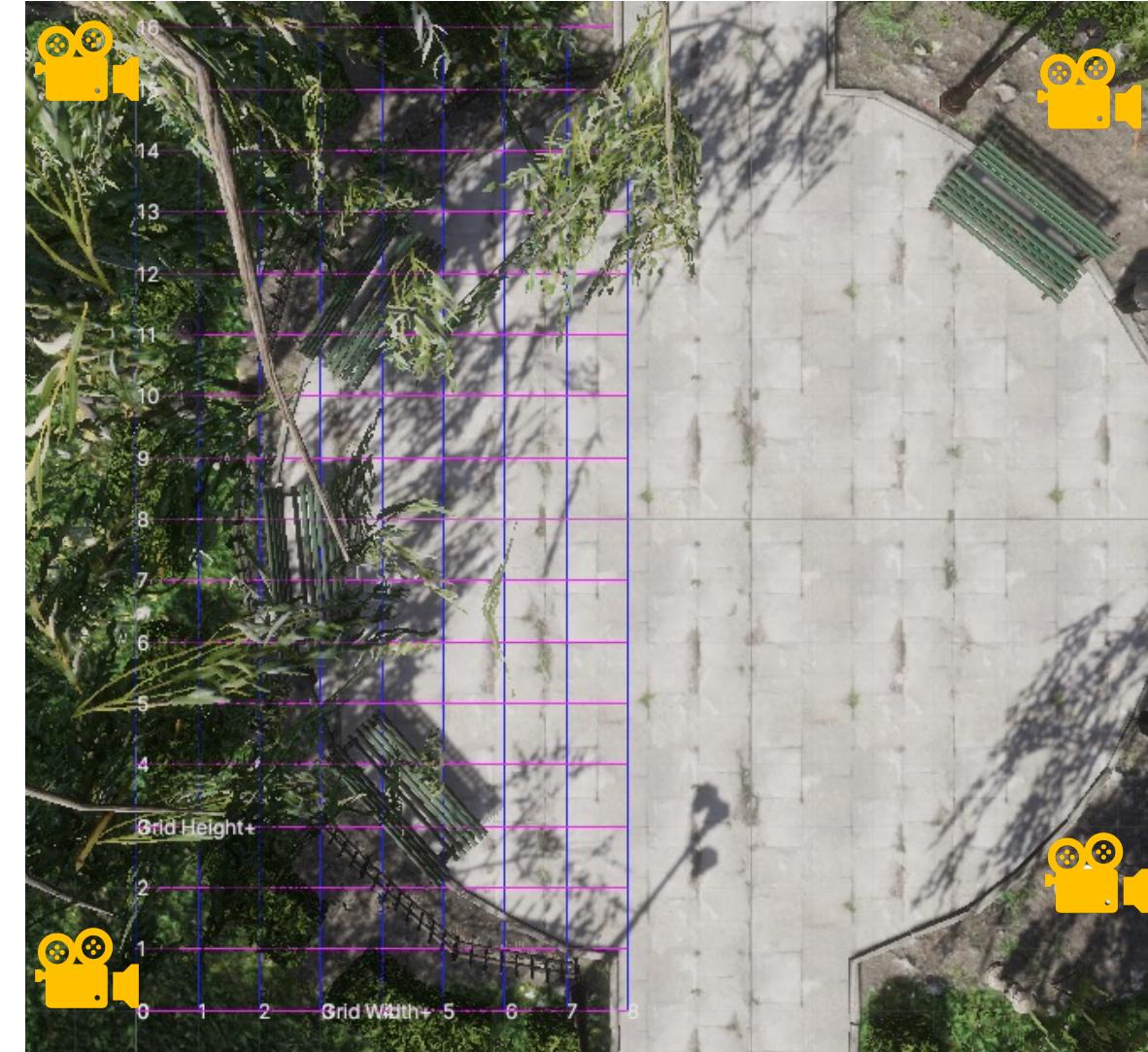
Ground Plane AoI (Area of Interest)

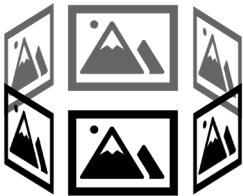
Grid resolution
0.025 m

Ground plane height = 15 m

Grid height = 600

Ground plane width = 10 m
Grid width = 400



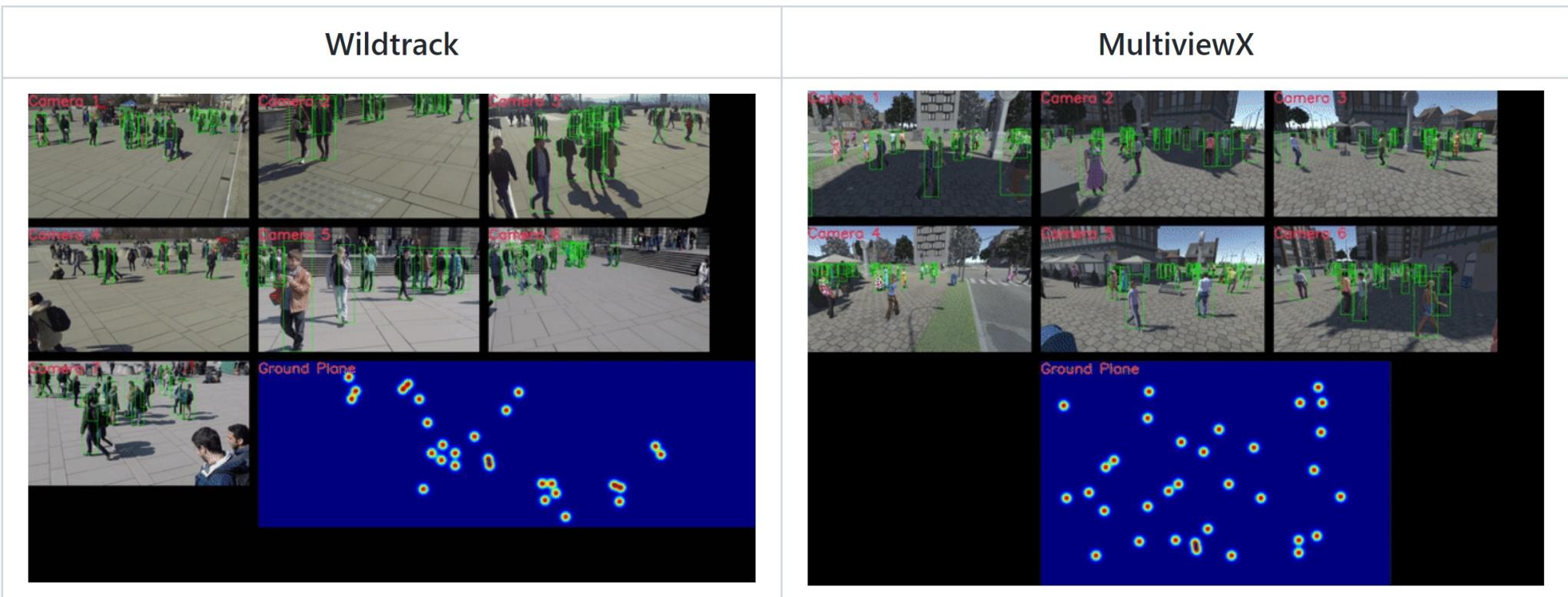


Preliminaries



Multi-view Datasets

	#camera	resolution	frames	area	crowdedness	avg. coverage
Wildtrack	7	1080×1920	400	$12 \times 36 \text{ m}^2$	20 person/frame	3.74 cameras
MultiviewX	6	1080×1920	400	$16 \times 25 \text{ m}^2$	40 person/frame	4.41 cameras





Preliminaries



Metrics (2D Localization)

$$MODA = 1 - \frac{FP + FN}{N} \quad \rightarrow \quad \text{Considers both false positives and false negatives.}$$

$$MODP = \frac{\sum 1 - \frac{d[d < t]}{t}}{TP} \quad \rightarrow \quad \text{Considers localization precision of true positives.}$$

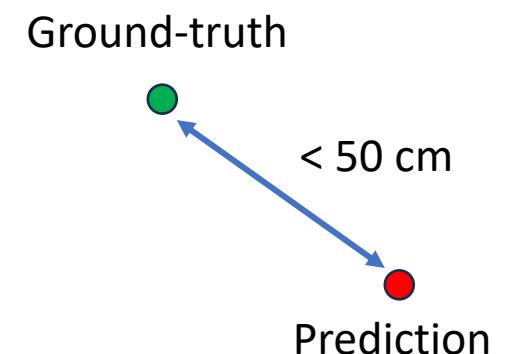
$$Precision = \frac{TP}{FP + TP} \quad \rightarrow \quad \text{Considers true positive rate.}$$

$$Recall = \frac{TP}{N} \quad \rightarrow \quad \text{Considers accurate localization performance.}$$

d = distance from a detection to its ground truth

t = threshold (50 cm)

N = total number of detections in the ground truth

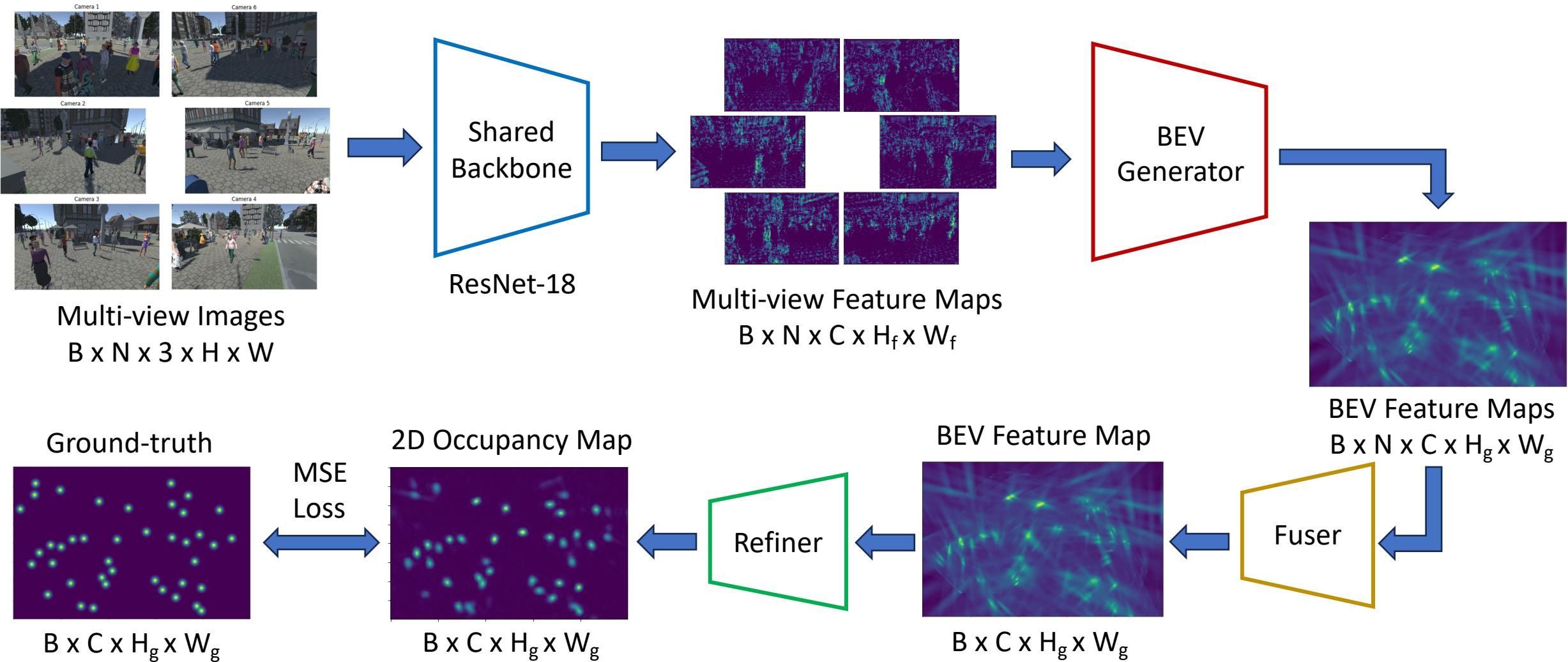




Preliminaries



Overall Framework





Preliminaries



Voxelization

Voxel Size = 0.01

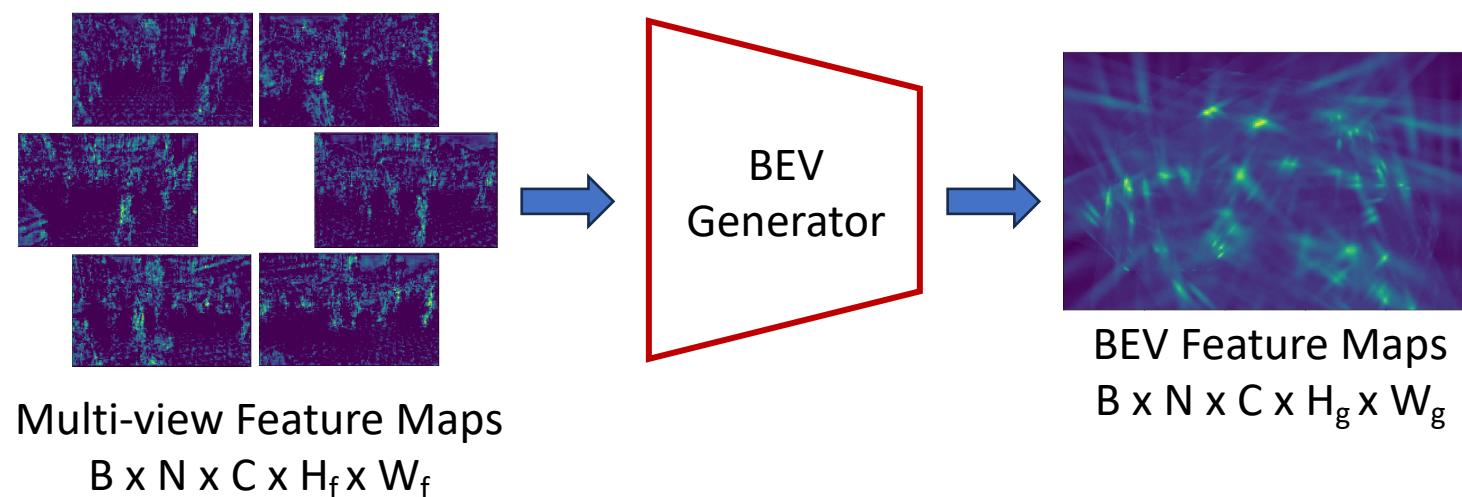




Related Work



2D to BEV Feature Transform

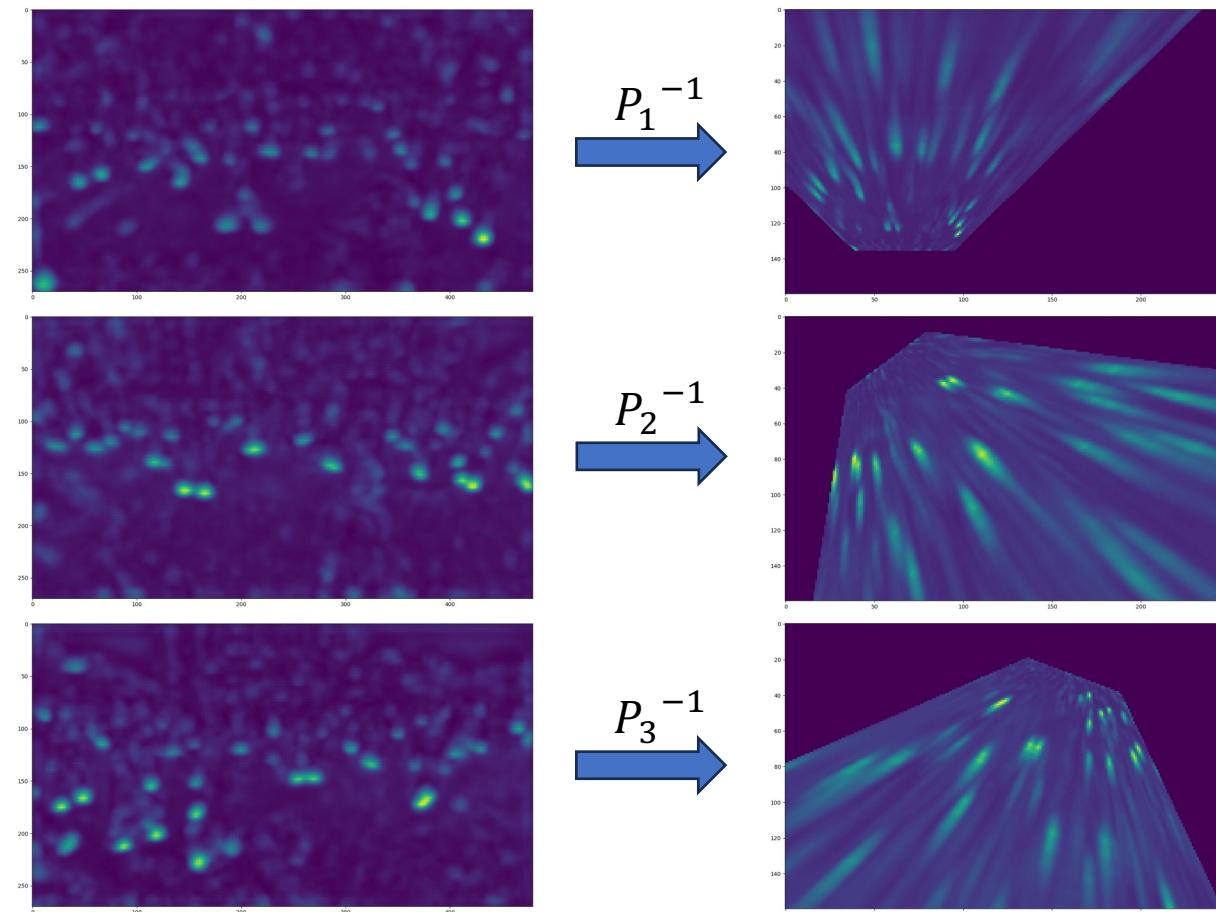




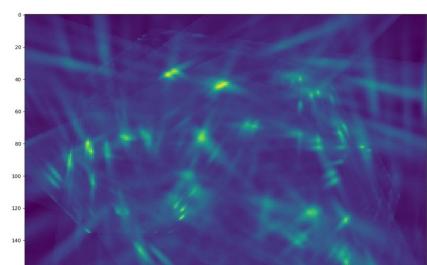
Related Work



Inverse Perspective Mapping (MVDet ECCV 2020)



Aggregated BEV Feature



$B \times NC \times H_g \times W_g$

Concatenate

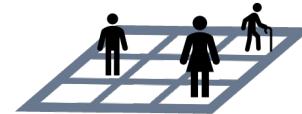
$$z \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R \ T] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

$$P = K[R \ T]$$

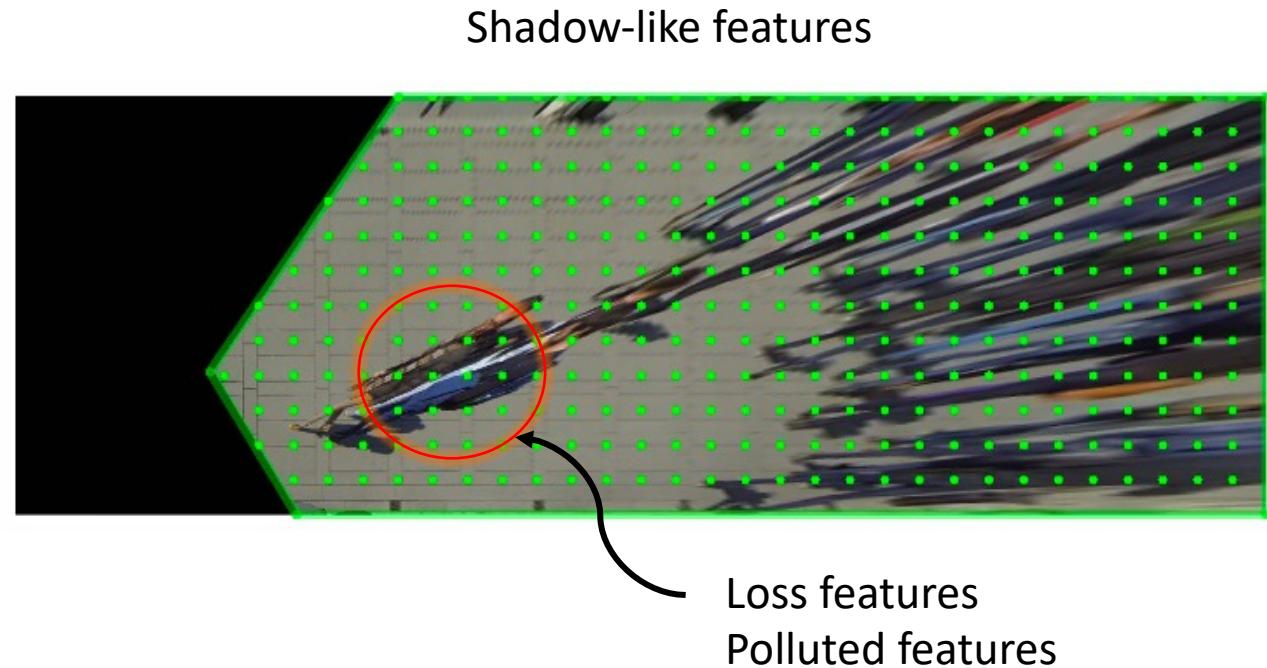
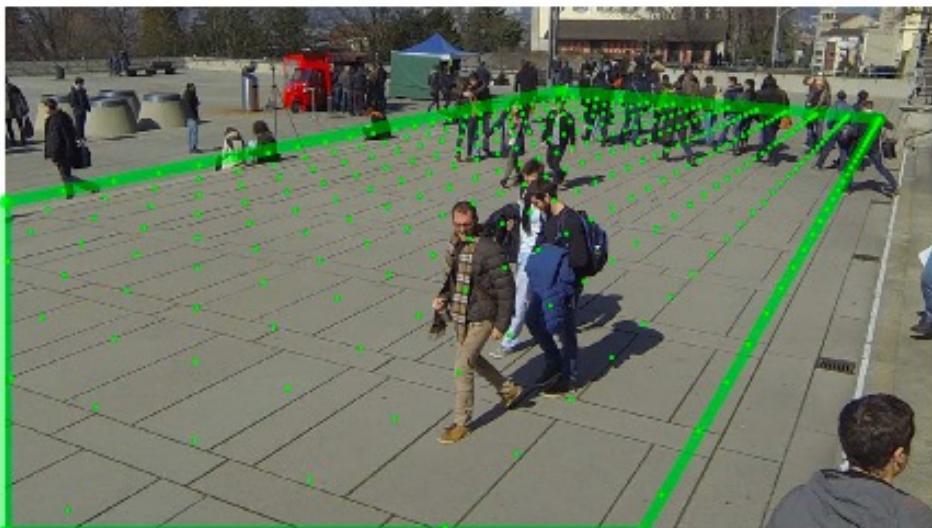
$$X_w = P^{-1} X_i$$



Related Work



Inverse Perspective Mapping (MVDet ECCV 2020)

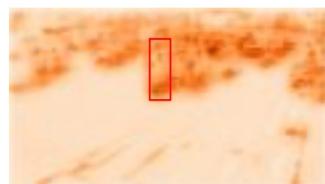




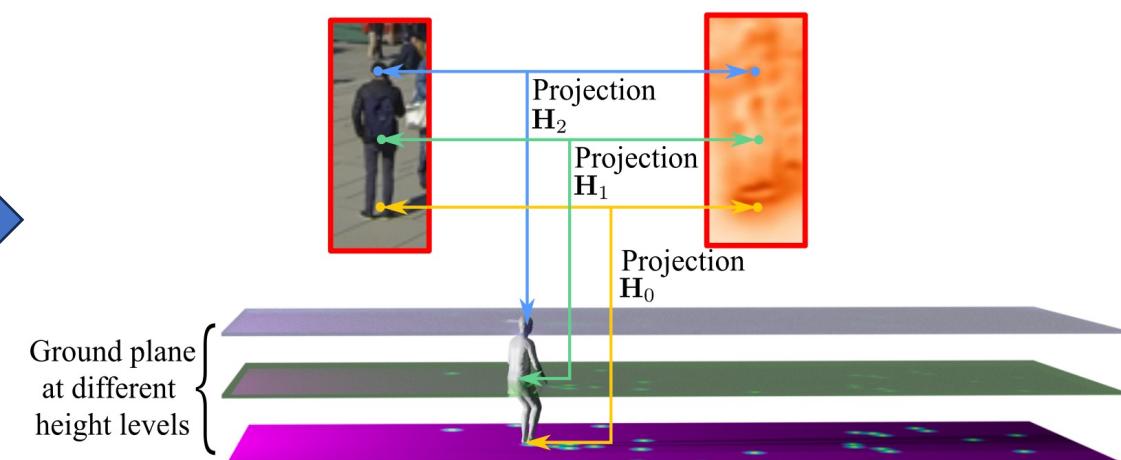
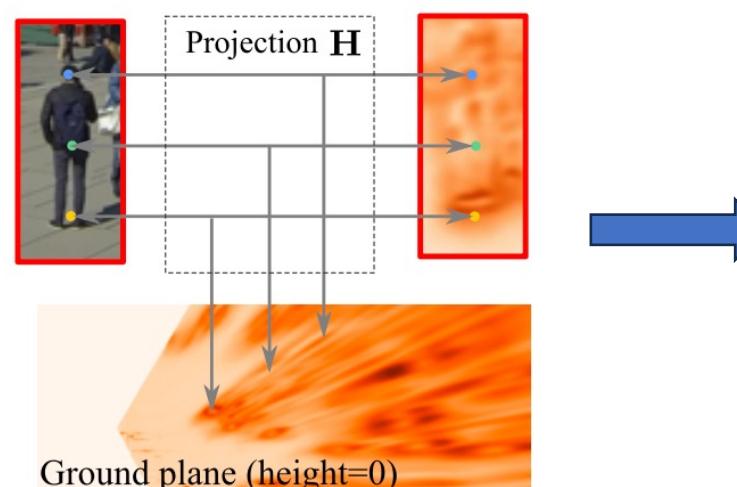
Related Work



Inverse Perspective Mapping with Multiple Homographies (SHOT ICCV 2021)



Feature Map



$$z \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R \quad T] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

0m
0.9m
1.8m



Related Work



3D Feature Lifting

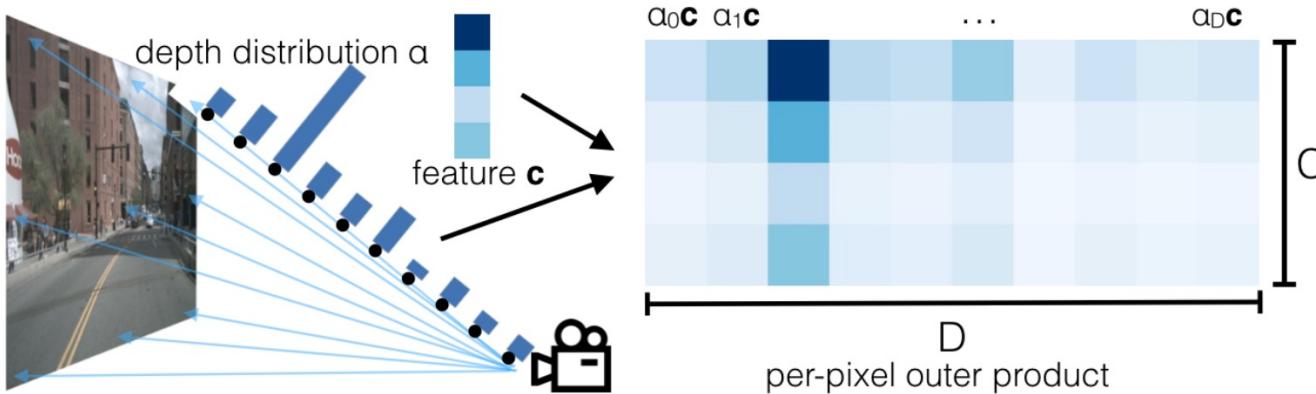
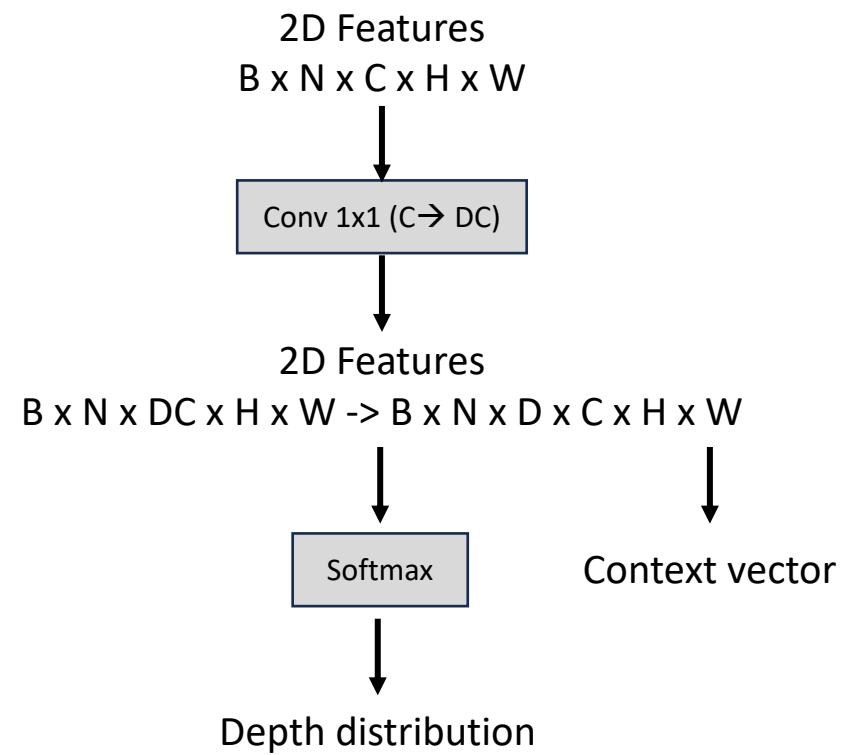
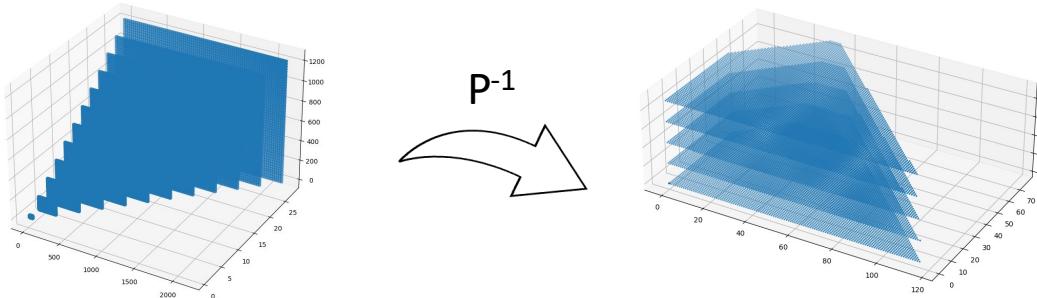


Fig. 3: We visualize the “lift” step of our model. For each pixel, we predict a categorical distribution over depth $\alpha \in \Delta^{D-1}$ (left) and a context vector $\mathbf{c} \in \mathbb{R}^C$ (top left). Features at each point along the ray are determined by the outer product of α and \mathbf{c} (right).



$$z \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R \quad T] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$



Related Work



3D Feature Pulling

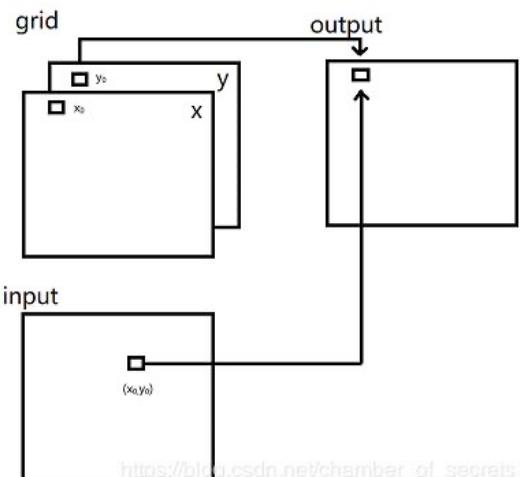
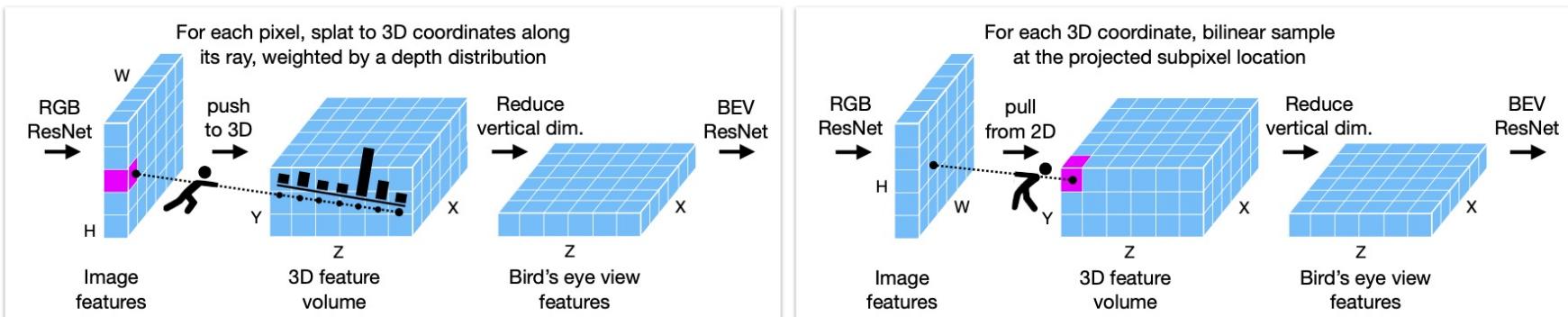
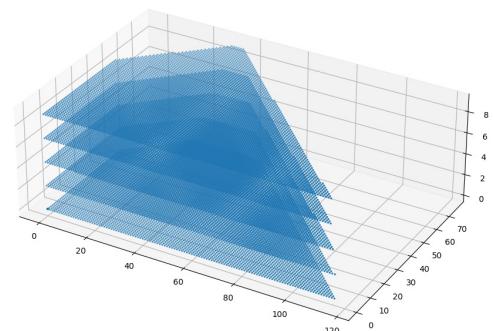
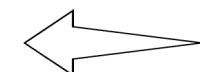
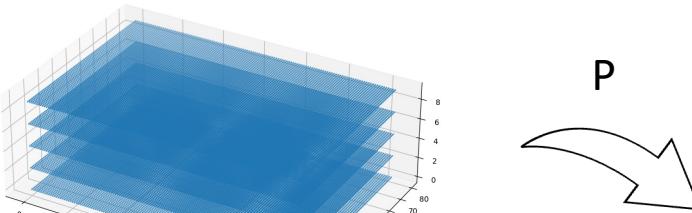


Fig. 1: **2D-to-BEV architecture, illustrated with two lifting strategies.** The left panel shows the Lift-Splat approach [9]: in this method, each 2D feature is “pushed” to 3D, filling voxels that intersect with its ray. The right panel shows our bilinear sampling approach: in this method, each 3D voxel “pulls” a feature from the 2D map, by projection and subpixel sampling.

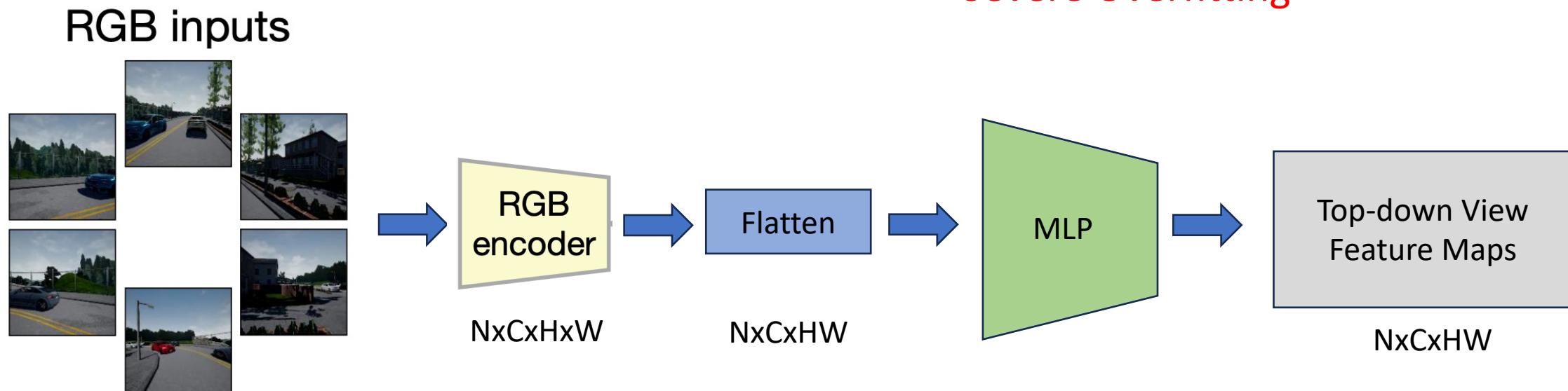




Related Work

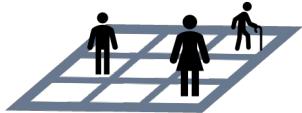


MLP-based Un-projection

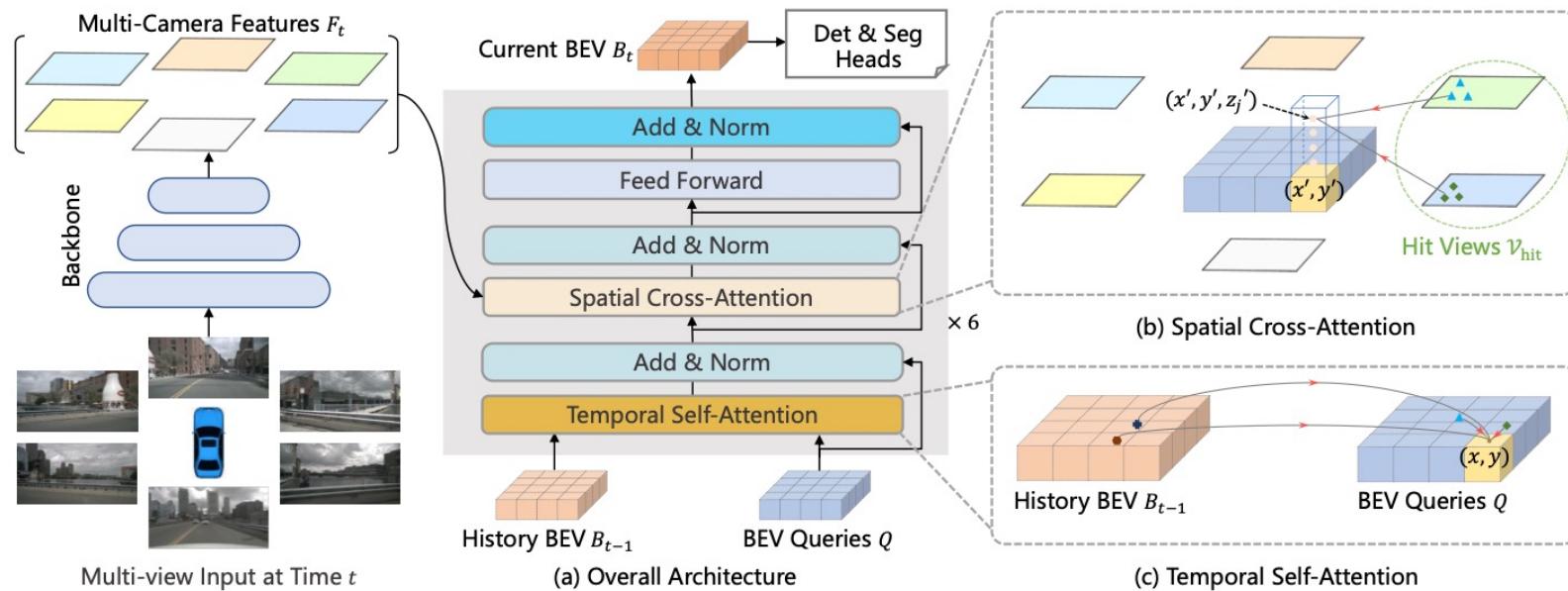




Related Work



Geometry-aware Transformer-like Models (Explicit-based)



$$\text{DeformAttn}(q, p, x) = \sum_{i=1}^{N_{\text{head}}} \mathcal{W}_i \sum_{j=1}^{N_{\text{key}}} \mathcal{A}_{ij} \cdot \mathcal{W}'_i x(p + \Delta p_{ij}),$$

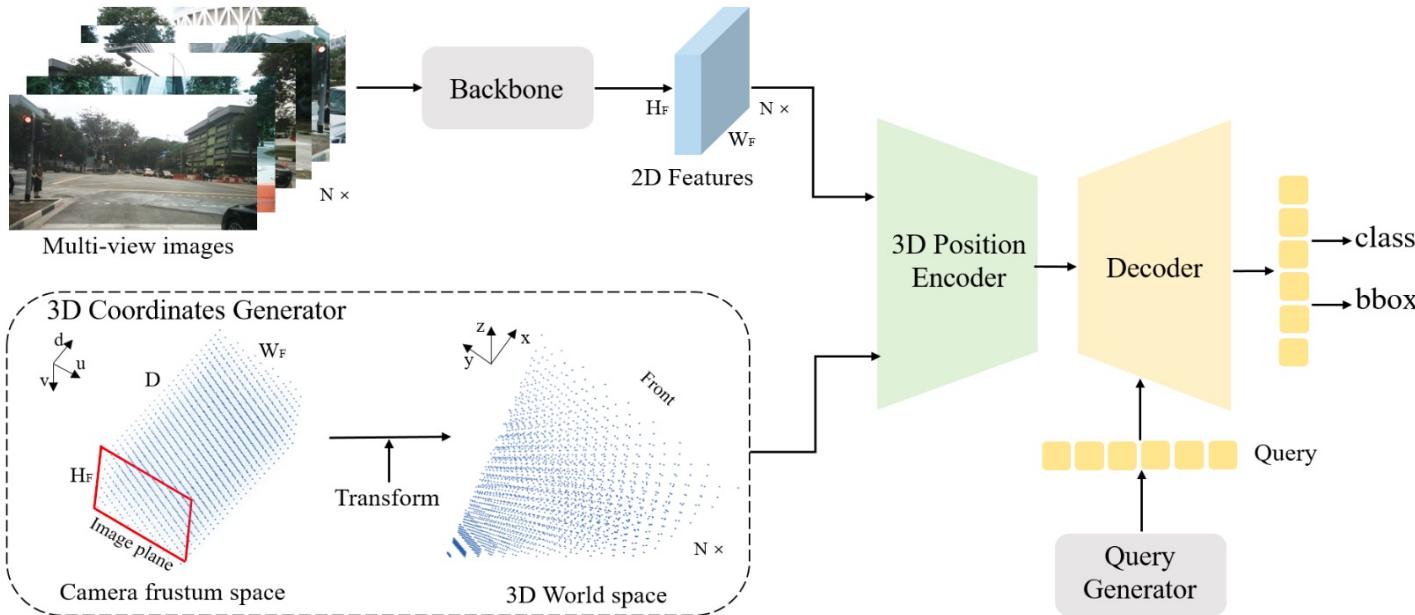
$$\text{SCA}(Q_p, F_t) = \frac{1}{|\mathcal{V}_{hit}|} \sum_{i \in \mathcal{V}_{hit}} \sum_{j=1}^{N_{\text{ref}}} \text{DeformAttn}(Q_p, \mathcal{P}(p, i, j), F_t^i),$$



Related Work

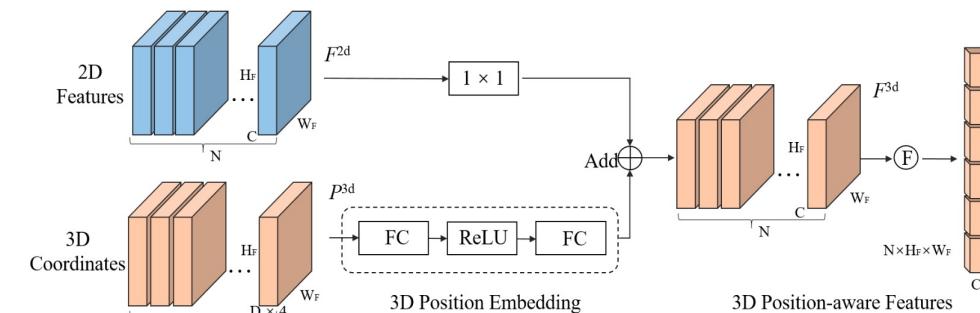


Geometry-aware Transformer-like Models (Implicit-based)

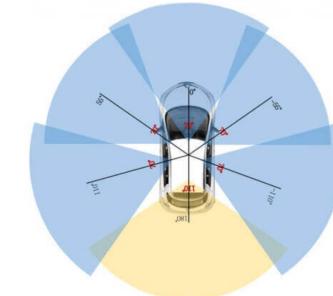


$$p_{i,j}^{3d} = K_i^{-1} p_j^m \quad P^{3d} = \{P_i^{3d} \in R^{(D \times 4) \times H_F \times W_F}, i = 1, 2, \dots, N\}$$

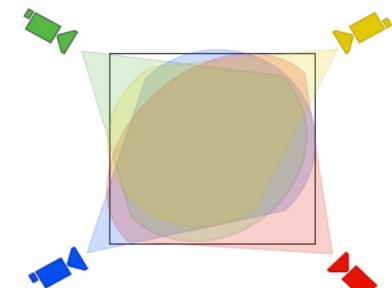
$$F_i^{3d} = \psi(F_i^{2d}, P_i^{3d})$$



Autonomous Driving



Our case





Related Work



2D to BEV Feature Transformation Summary

- Learnable
 - MLP
 - Deformable attention
 - 3D Feature Lifting
- Parameter-free
 - Inverse Perspective Mapping
 - 3D Feature Pulling

TABLE I: Effect of lifting strategy.

2D-to-BEV strategy	IOU
Unweighted splatting	43.1
Depth-based splatting [9]	44.4
Deformable attention [5]	46.5
Bilinear sampling (ours)	<u>47.4</u>
Multi-scale deform. attn. [5]	48.9



Previous Work



Requirements for Multi-view Detection System

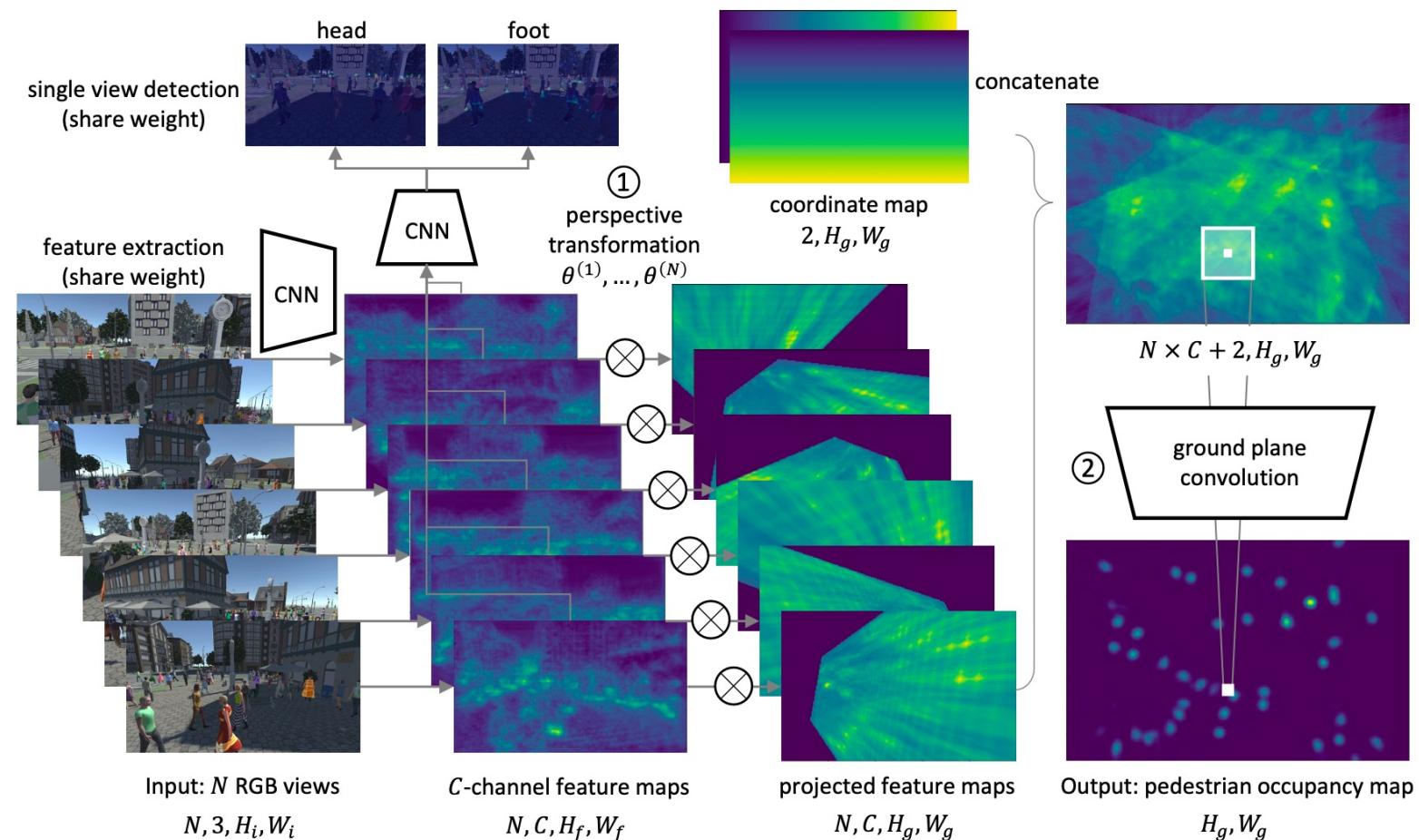
- Dataset size is small (400 frames) -> Overfitting -> Model size needs to be small
- Crowded scene -> High resolution feature map
- Generalization performance



Previous Work



MVDet (ECCV 2020)

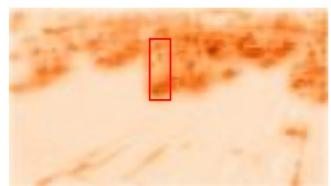




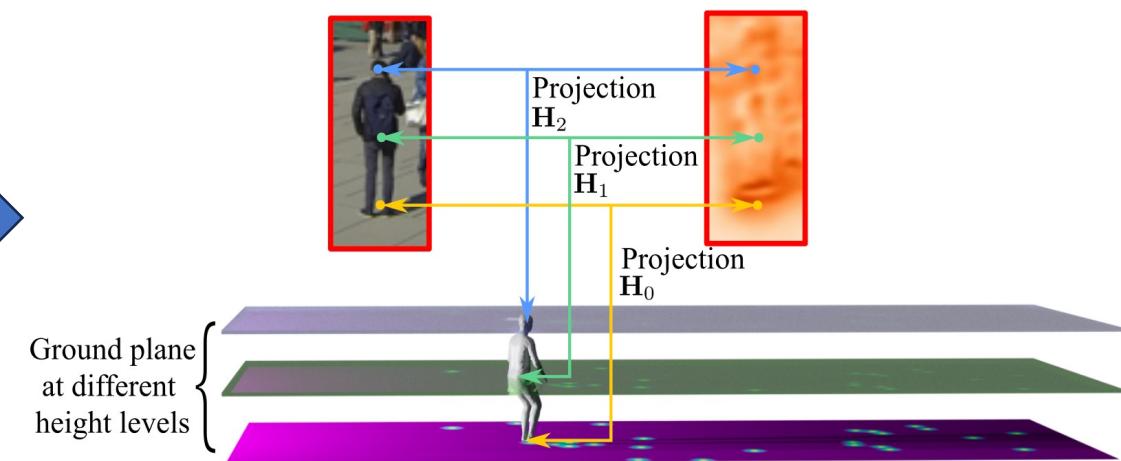
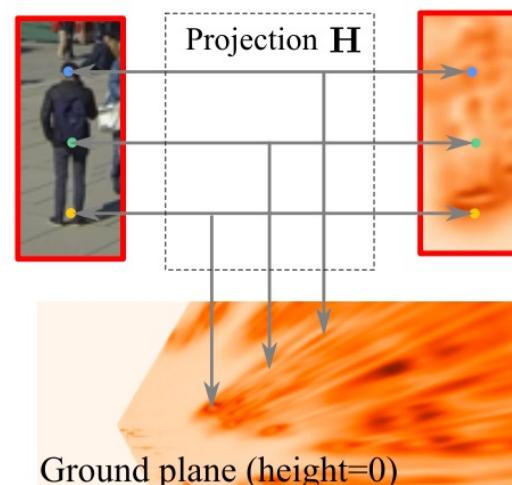
Previous Work



SHOT (ICCV 2021)

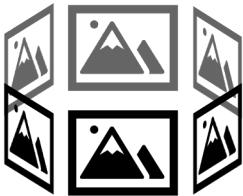


Feature Map

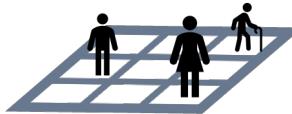


$$z \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R \quad T] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

	MODA	MODP	P	R
MVDet	83.9	79.6	96.8	86.7
SHOT	88.3 (+4.4)	82.0	96.6	91.5



Previous Work



3DROM (ECCV 2022)

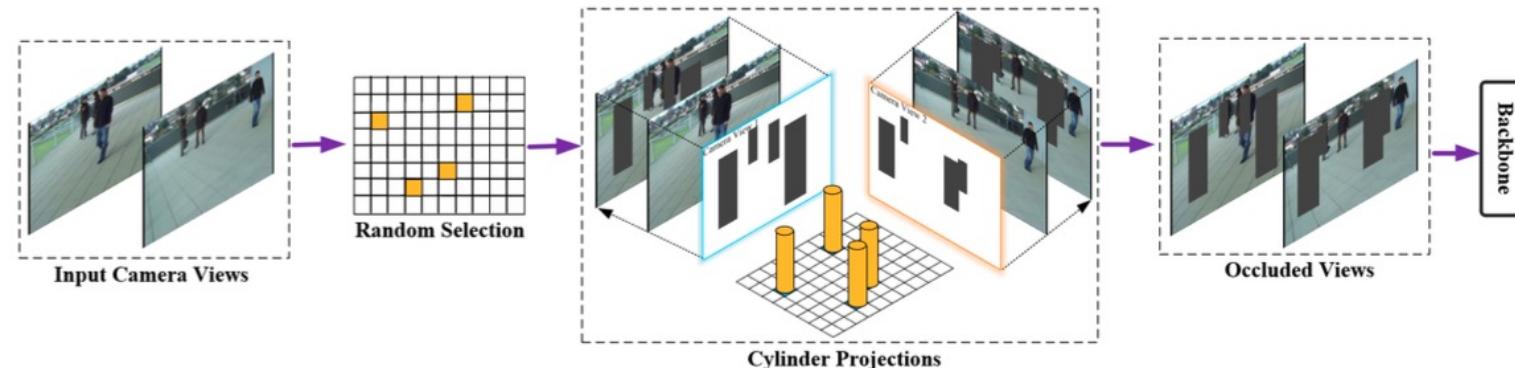
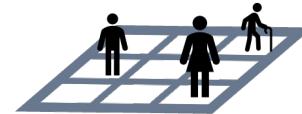


Fig. 2. A schematic diagram of the 3D Random Occlusion method.

	MODA	MODP	P	R
MVDet	83.9	79.6	96.8	86.7
SHOT (Multiple Homographies)	88.3 (+4.4)	82.0	96.6	91.5
3DROM (SHOT + 3D Random Erasing)	95.0 (+6.7)	84.9	99.0	96.1



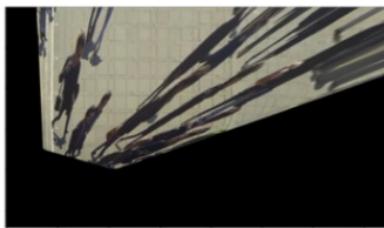
Previous Work



MVAug (WACV 2023)



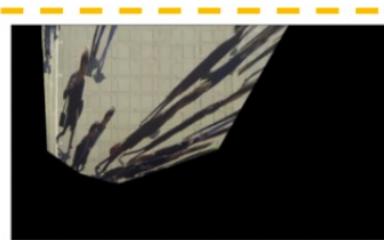
Original image



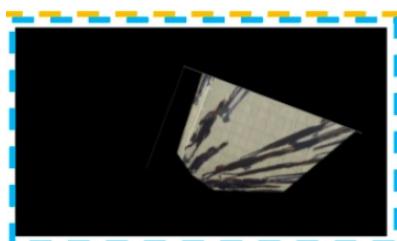
Original projection



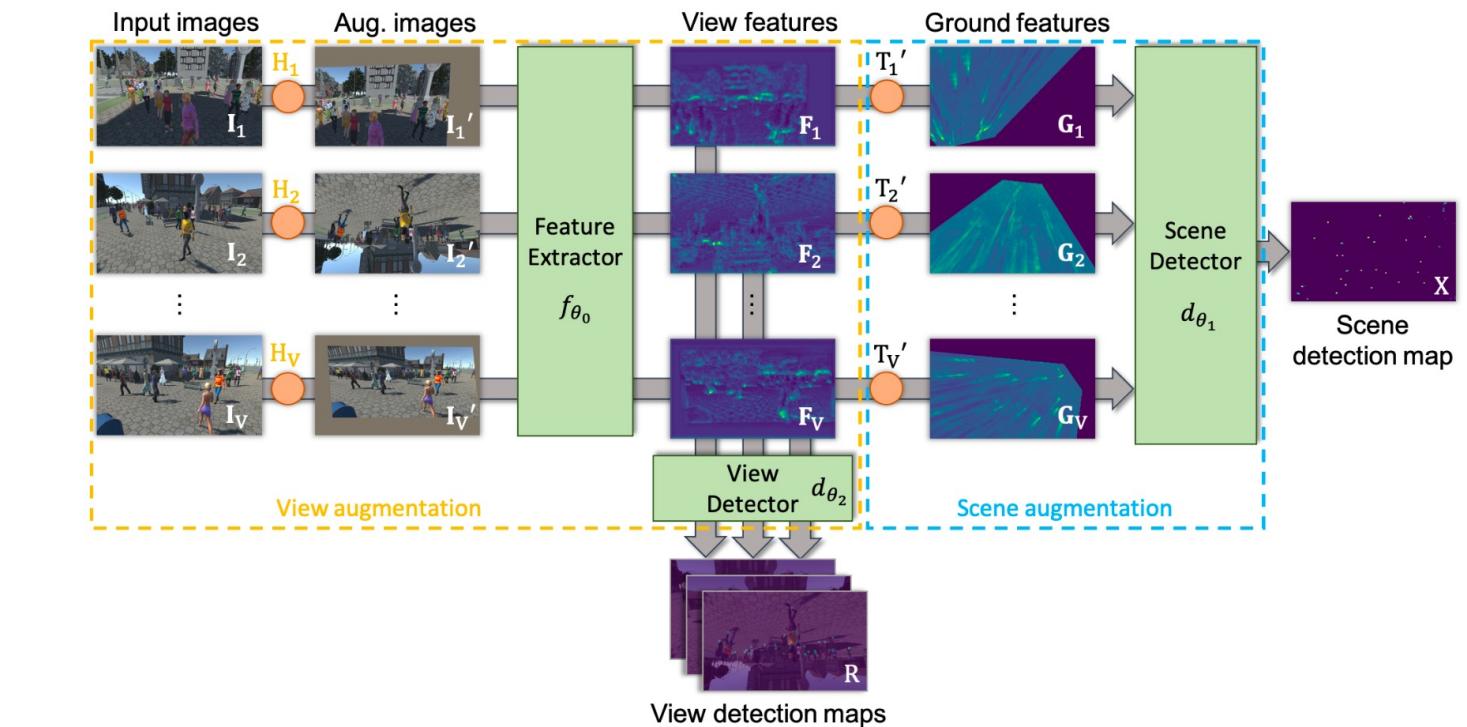
View aug.



View aug. projection



Both aug. projection



	MODA	MODP	P	R
MVDet	83.9	79.6	96.8	86.7
MVAug	95.3 (+11.4)	89.7	99.4	95.9

