

Multiview Detection with Shadow Transformer (and View-Coherent Data Augmentation)

Yunzhong Hou

Australian National University
Canberra, ACT 2600, Australia
yunzhong.hou@anu.edu.au

Liang Zheng*

Australian National University
Canberra, ACT 2600, Australia
liang.zheng@anu.edu.au

ABSTRACT

Multiview detection incorporates multiple camera views to deal with occlusions, and its central problem is multiview aggregation. Given feature map projections from multiple views onto a common ground plane, the state-of-the-art method addresses this problem via convolution, which applies the same calculation regardless of object locations. However, such translation-invariant behaviors might not be the best choice, as object features undergo various projection distortions according to their positions and cameras. In this paper, we propose a novel multiview detector, MVDeTr, that adopts a newly introduced shadow transformer to aggregate multiview information. Unlike convolutions, shadow transformer attends differently at different positions and cameras to deal with various shadow-like distortions. We propose an effective training scheme that includes a new view-coherent data augmentation method, which applies random augmentations while maintaining multiview consistency. On two multiview detection benchmarks, we report new state-of-the-art accuracy with the proposed system. Code is available at <https://github.com/hou-yz/MVDeTr>.

CCS CONCEPTS

• Computing methodologies → Object detection.

KEYWORDS

multiview detection; transformer; data augmentation

1 INTRODUCTION

Incorporating multiple camera views effectively alleviates occlusion [1, 14, 26], a challenging problem in the community. For monocular view systems, occlusion prevents the system from discovering targets hidden behind obstacles or other people. On the other hand, multiview systems usually leverage multiview images that are synchronized and have overlapping fields of view, and thus complement each other in inferring what is occluded. Camera calibrations are also presumably provided in multiview detection, which enable aggregation from multiple views on the ground plane (a common reference plane) of the bird’s eye view.

A central problem in multiview detection is how to aggregate content from multiple camera views. A recent state-of-the-art method, MVDet [14], adopts a fully-convolutional approach. It first projects feature maps onto the ground plane and then applies convolution to jointly consider neighboring locations over multiple camera views. The convolution-based fusion approach, being translation-invariant (same calculations regardless of locations and cameras), arguably does not well suit the projected feature maps. In fact, different

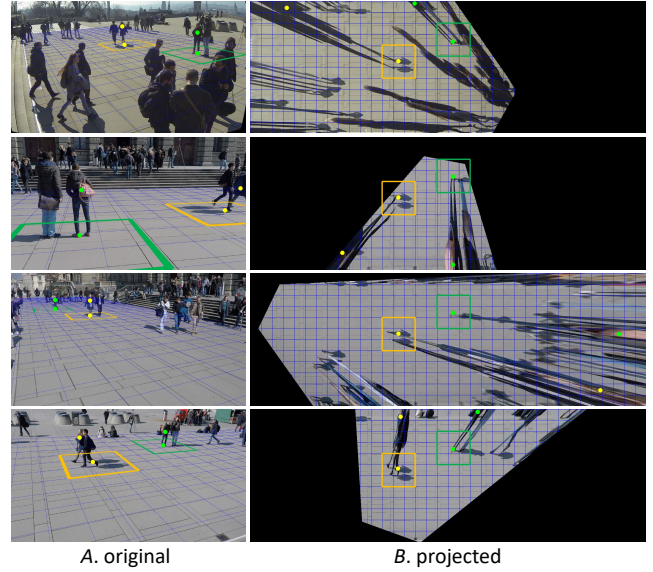


Figure 1: Distinct distortion patterns at different locations vs. fixed computation in convolution. We highlight two pedestrians with yellow and green dots as example. *A.* original images from multiple cameras. *B.* projected images onto the ground plane with bird’s eye view, where we observe shadow-like distortion patterns that vary with the locations and views. Currently, convolution (yellow and green boxes in *B* and their corresponding receptive fields in *A*) is adopted to aggregate across multiple projected feature maps (we use RGB images here for better illustration). Due to its translation invariance, convolution applies the same calculation at all locations, even when distortion patterns at different positions might be rather different. Therefore, we investigate other design choices for multiview aggregation, which can attend adaptively at different locations, thus better accommodating the shadow-like distortion patterns.

positions and different cameras in the scenario might introduce diverse distortion patterns during projection, just as different locations and different street lights cause pedestrians to cast various shadows. As shown in Fig. 1, pedestrians (and their features) are projected rather dis-similarly across different positions and cameras. Therefore, applying the same operation with convolution (due to its innate characteristic of *translation invariance*) potentially contradicts the distinct *position-sensitive* distortion patterns, further limiting the system performance.

*Corresponding author

In this paper, we propose MVDeTr, an end-to-end multiview detector that deals with various distortion patterns during multiview aggregation using shadow transformer. Originated from the deformable transformer [47], shadow transformer is designed to attend to different pixels across multiple cameras from a certain position and view. Departing from deformable transformers that attend within a single image, shadow transformer jointly considers multiple projected feature maps when updating features of each view. Specifically, we feed the projected feature maps, a shared ground plane position embedding, and the learned camera embeddings to the deformable transformer, which make the system aware of the object, position, and camera characteristics all at once. This view- and position-aware attention mechanism can better address projection distortions (*e.g.*, skews and stretches) rarely observed in natural images. Shadow transformer only has an encoder component, and the full system (MVDeTr) adopts a fully connected layer to output the detection result for each ground plane location.

Together with the proposed architecture, a new training scheme is introduced that includes a new view-coherent data augmentation approach and loss terms from related monocular-view detection systems [16, 19, 46]. Despite its success in preventing overfitting and allowing higher accuracy, data augmentation is rarely applied in multiview detection systems. This is because of the restriction that the augmented input images must maintain multiview consistency: for joint consideration of multiple camera views to benefit the detection system, these views must be coherent in the content they provide. Since the augmentation policies are mostly random, the augmented input views are likely to violate the consistency restriction, rendering them unusable in multiview systems. In the proposed view-coherent data augmentation, we use different random affine transformations for different views, and extract image features from these augmented images. We also adjust the ground truth accordingly for per-view supervision. Then, we invert the affine transformations before projecting the feature maps onto the ground plane, thus restoring the multiview coherency for ground plane supervision when applying random data augmentations.

On Wildtrack [4], a real-world dataset, and MultiviewX [14], a synthetic dataset, MVDeTr achieves state-of-the-art detection accuracy. Specifically, MVDeTr requires a similar number of training iterations and memory footprints compared to the convolution-based alternatives.

2 RELATED WORK

Monocular-view detection. The majority of works in object detection focus on the monocular-view setting, where only a single image is available. Many inspiring methods are proposed for this task, including two-stage detectors like the R-CNN family [11, 24] that regress bounding boxes from proposals, and one-stage detectors that adopt the anchor-based [19] or anchor-free [33, 46] approach. Some also investigated other detection like point cloud [23] and more [18, 27, 30, 39, 43]. In addition to the architectures, many novel training schemes are also investigated. To list a few, RezaTofighi *et al.* [25] propose a novel localization loss, G-IOU, directly derived from intersection-over-union for object detector training. Lin *et al.* [19] propose Focal Loss to deal with the foreground-background imbalance in single-stage detectors. Many researchers focus on the

pedestrian detection problem, rather than general-case object detection. To deal with occlusions, some investigate part-based methods, as the occluded human are still partially observable [21, 22, 32, 44], while others investigate learned NMS [13], specific loss design like Repulsion Loss [37], and more [17, 38, 45].

Multiview detection. In order to detect pedestrians under heavy occlusion, many researchers choose to explore multiple camera views. These views are synchronized and calibrated, providing complementary descriptions of the scene. Note that camera calibrations also produce a correspondence between each location on the ground plane and its bounding boxes in multiple camera views, as we can calculate the 2D bounding boxes once we assume an average 3D human height and width via perspective transformation. In this case, multiview detection systems usually evaluate their performance with pedestrian occupancy maps on the ground plane [1, 4, 14]. To aggregate from multiple views, the core problem in multiview detection, researchers adopt different approaches. Baque *et al.* [1] use higher-order potentials to model the consistency between spatially neighboring locations, but requires additional potential terms design and calculation outside neural networks. Hou *et al.* [14] also considers the spatially neighboring locations in multiview aggregation, but accomplish this via convolution and achieve state-of-the-art performance. Of all these methods, multiview aggregation either needs specific design outside neural network [1, 10], or fail to consider spatial neighboring locations [3], or adopt computations un-adjustable to the locations [14]. With that said, none of the current systems applies data augmentations in their training scheme, mainly due to the restriction of multiview coherency.

Transformers and their application. Recently, researchers investigate transformers [7, 35] and their applications [2, 8, 31, 41, 42]. Originally applied in sequence reasoning tasks, transformers adopt multi-head self-attention to model the relationship between each point. Compared to Non-local networks [36], it uses fully connected layers to improve model capability, while providing location sensitivity by feeding position embedding. On natural language tasks including machine translation and question answering, transformers achieve great performance, proving their scalability and capability. For vision tasks, Dosovitskiy *et al.* propose ViT as a transformer-based image classifier [8]. In object detection, DETR [2] use a transformer to model the object relationship in feature maps, and adopt a novel end-to-end object detection training scheme departing from current ones. Given the long training time and large computation overhead of DETR, Zhu *et al.* propose Deformable DETR, which reduces the computation by only attending to several points around the reference rather than the entire feature map. Similar to Deformable convolution [5], the attended points in Deformable DETR are also fully learnable.

3 METHOD

In this section, we first revisit the multiview aggregation problem and deformable transformer. Then, we propose the MVDeTr architecture that enables adaptive attention across different locations and cameras. Lastly, we introduce the new training scheme for MVDeTr that allows random data augmentation in multiview detection while maintaining multiview coherency.

3.1 Revisiting Multiview Aggregation and Deformable Transformer

Multiview aggregation with feature projection. In the recent state-of-the-art, Hou *et al.* [14] adopt an anchor free representation, by projecting the multiview feature maps to the ground plane (bird’s eye view). Given 2D image pixel coordinate (u, v) , we calculate the 3D world position $(x, y, z)|_{z=0}$ on the ground plane ($z = 0$) with

$$\gamma \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = P_\theta \begin{pmatrix} x \\ y \\ z \end{pmatrix} = A [R|t] \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} & \theta_{14} \\ \theta_{21} & \theta_{22} & \theta_{23} & \theta_{24} \\ \theta_{31} & \theta_{32} & \theta_{33} & \theta_{34} \end{bmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}, \quad (1)$$

where γ is a scaling factor, and P_θ is the perspective transformation matrix calculated using the intrinsic parameter A and rotation-translation matrix (extrinsic parameter) $[R|t]$. By setting $z = 0$, we can retrieve the correspondence between image pixel (u, v) and the ground plane position (x, y) , allowing for feature map projection.

By concatenating all projected feature maps, the system then gains access to a feature description for each location across all cameras. This anchor-free representation outperforms previous anchor-based ones [1, 3] that rely on ROI-pooling [24] over the anchor boxes, as anchor boxes might not be accurate in the first place and can lead to inaccurate feature representations.

Moreover, researchers find that jointly considering neighboring locations with conditional random field (CRF) [1] or convolution [14] also helps with the multiview aggregation. But these methods also have their drawbacks: CRF requires specific potential terms design and operations outside neural network forward pass, whereas convolution is translation-invariant and does not suit the various distortion patterns from projection (see Fig. 1).

Deformable transformer. Proposed in Deformable DETR [47], deformable transformer departs from previous self-attention-based transformer architectures [35] that attend to all locations, and only focuses on a few selected reference points. Specifically, given feature map f and a position p , for each head $m \in \{1, \dots, M\}$ in the M -head pipeline, deformable attention first outputs position offsets $\{\Delta p_{mk}\}_{k=1}^K$ for K reference points. Then, it estimates attention weights $\{a_{mk}\}$ for the reference points $\{\Delta p_{mk}\}$ in each head ($\sum_k a_{mk} = 1$), and outputs the final results over multi-head,

$$\text{DeformAttn}(f, p) = \sum_{m=1}^M W_m \sum_{k=1}^K a_{mk} W'_m f(p + \Delta p_{mk}), \quad (2)$$

where W_m and W'_m denote the translations for head m .

Such deformable multi-head attention allows for adaptive sampling positions $\{p + \Delta p_{mk}\}$ for each position p . However, it still lacks specific designs so as to aggregate from multiple camera views rather than a single image.

3.2 MVDeTr with Shadow Transformer

In this paper, we propose MVDeTr, a novel multiview detector that adopts a so-called shadow transformer to aggregate from multiple projection-distorted feature maps (Fig. 1). Rather than a single image as usually does in deformable transformers, shadow transformer jointly considers multiple camera views. The overall architecture of MVDeTr is shown in Fig. 2.

Given inputs from multiple camera views, first, MVDeTr extracts image feature maps via a ResNet-18 [12] network. Next, similar to the previous method [14], we project the feature maps onto the ground plane to enable further multiview fusion. Then, cues from multiple views are jointly and adaptively processed with the shadow transformer, where a multiview-deformable-attention-based transformer encoder is adopted. We do not use any transformer decoder as we still formulate the multiview detection problem as a key point regression problem [16, 46] rather than as a direct set prediction problem [2] (see Section 3.3 for more details). Lastly, outputs from multiple views are aggregated using a fully connected layer. It is noteworthy that the joint consideration of multiple views only relies on the shadow transformer and the last fully connected layer. In this manner, MVDeTr fuses multiview features with *position-sensitive* calculations, in contrast to the *translation-invariant* convolution pipeline that applies the same calculation to all positions.

Specifically, multiview deformable attention is the core of the proposed shadow transformer. Inspired by multi-scale deformable attention [47], we extend deformable attention to multiview detection by jointly considering multiple projected feature maps. For C camera views, multiview deformable attention takes as input the projected feature maps $\{f_c\}_{c=1}^C$ with shadow-like distortions, a certain position p (total number of p is the size of the ground plane feature map) across all C feature maps, and learned camera embeddings for each camera. Then, it translates each view into M heads with view-shared fully connected layers, and jointly considers all $M \times K \times C$ reference points when updating for each position at each view,

$$\begin{aligned} & \text{MVDeformAttn}(\{f_c\}_{c=1}^C, p, c) \\ &= \sum_{m=1}^M W_m \sum_{c'=1}^C \sum_{k=1}^K a_{mkc'} W'_m f(p + \Delta p_{mkc'}). \end{aligned} \quad (3)$$

Specifically, the attention weights satisfy ($\sum_{k,c} a_{mkc} = 1$). We use the multiview deformable attention to construct an encoder for multiview processing and aggregation in MVDeTr.

As shown in Fig. 2, the multi-head architecture allow the system to focus on different aspects. For instance, the blue branch focuses on the real shadows cast by the sun, thus attending similar relative positions across multiple views (human shadow cast by the sun are already on the ground plane, and further projection with Eq. 1 does not bring forward any distortion as shown in Fig. 1). On the other hand, the yellow branch focuses on the various “shadow-like” distortion patterns of human bodies and features, which are caused by perspective transformations at different positions from different cameras (human bodies are projected onto the ground plane, as if there is a light source at the camera location casting shadows onto the ground plane, thus giving distinct distortion patterns at different locations). Accordingly, the yellow branch attends to different points for different locations.

3.3 Training Scheme

Loss functions. We formulate multiview detection as a key point detection problem, as the multiview systems evaluate pedestrian occupancy on the ground plane [4, 14]. Following CornerNet [16], we regress a heatmap of the pedestrian occupancy likelihood score

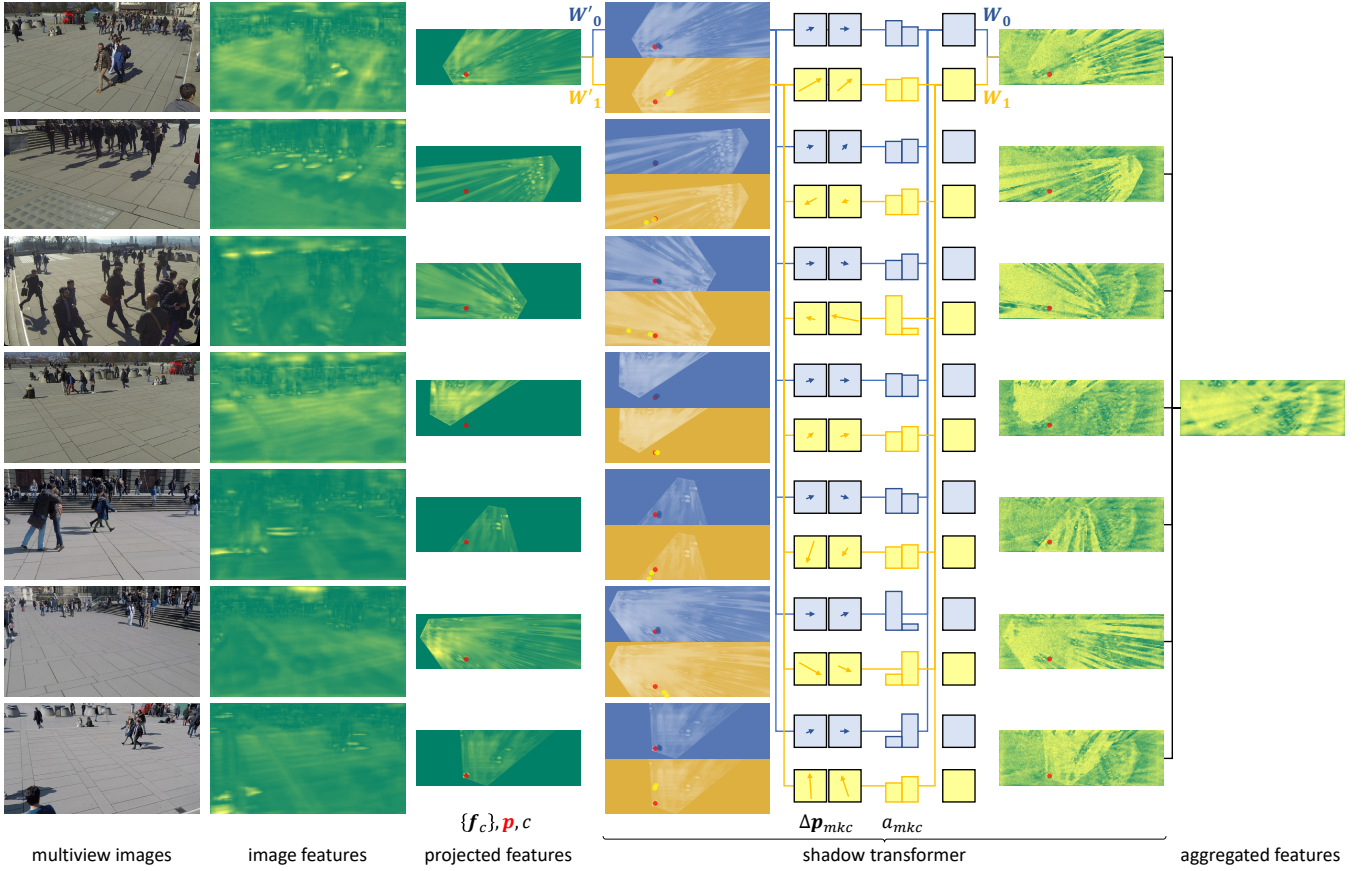


Figure 2: Multiview aggregation in MVDeTr. Given input images from C cameras, first, MVDeTr extracts image feature maps. Next, MVDeTr projects the feature maps onto the ground plane. Then, we adopt a shadow transformer to deal with the various shadow-like distortion patterns introduced by the perspective transformation. Last, we aggregate the output feature maps to formulate a final representation for the entire scenario. Specifically, inside this shadow transformer, we use multi-view deformable attention (an extension to the deformable attention [47]) to build an encoder for multiview processing and aggregation. To start with, we translate each of the projected feature maps $\{f_c\}$ into M heads. After that, at each position p , we estimate the position offsets of K reference points $\{p + \Delta p_{mkc}\}$ to attend on, and estimate their corresponding attention weights $\{a_{mkc}\}$ ($\sum_{k,c} a_{mkc} = 1$). Different from the original deformable transformer design, we consider all different cameras and M heads for each camera (which makes a total number of $M \times K \times C$ reference points for each location in each camera) for the multi-head attention, so as to infuse the multiview information to each of the views.

\hat{s}_p for each position p on the ground plane. Using Focal Loss [19] and a Gaussian smoothed target s ($s_p = 1$ for the ground truth location and gradually descent to 0), we write the detection loss as,

$$\mathcal{L}_{\text{det}} = \frac{1}{N} \sum_p \begin{cases} (1 - \hat{s}_p)^\alpha \log(\hat{s}_p) & \text{if } s_p = 1 \\ (1 - s_p)^\beta (\hat{s}_p)^\alpha \log(1 - \hat{s}_p) & \text{otherwise,} \end{cases} \quad (4)$$

where N denotes the total number of pedestrians over all positions. Besides, since the heatmap we output is often of lower resolution compared to the ground truth, we also regress an offset \hat{o}_p , which aims to compensate for the omitted decimal part $\frac{p}{r} - \lfloor \frac{p}{r} \rfloor$ during the $r \times$ downsampling. The offset loss can be written as,

$$\mathcal{L}_{\text{off}} = \frac{1}{N} \sum_{p|s_p=1} \left| o_p - \left(\frac{p}{r} - \left\lfloor \frac{p}{r} \right\rfloor \right) \right|, \quad (5)$$

where $|\cdot|$ denotes the \mathcal{L}_1 norm.

In addition to the detection loss \mathcal{L}_{det} and offset loss \mathcal{L}_{off} for the ground plane detection results, we also consider each of the individual images for supervision. Specifically, rather than detecting top-left and bottom-right corner [9, 16, 29] or the object center point [9, 20, 46], we detect the pedestrian foot point and regress its height and width. Foot points, or the center of the bounding box bottom, should be where the bounding box intersects with the ground plane ($z = 0$). Therefore, foot points on each image should be the same as or close to the corresponding pedestrian occupancy on the ground plane, and detecting foot points on each image should help to regress the ground plane occupancy. Additionally, estimating bounding box shapes also provides additional supervision, as the ground plane supervision does not include pedestrian height and

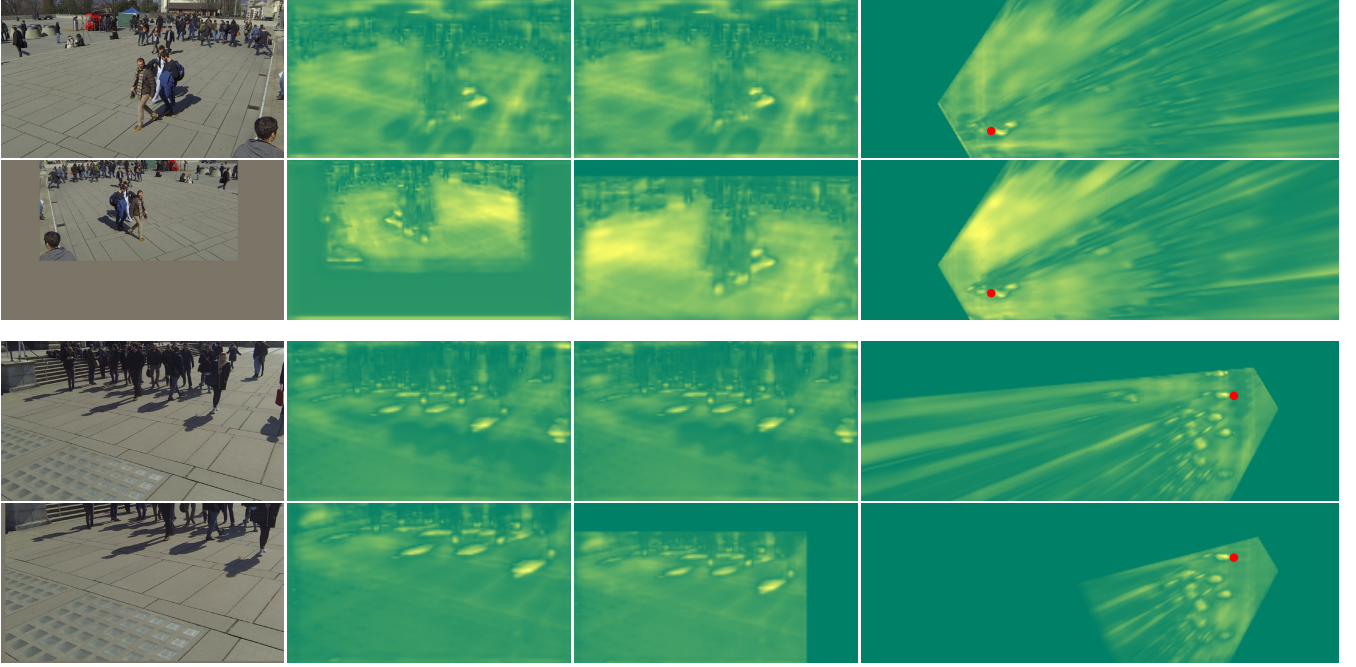


Figure 3: View-coherent data augmentation. In each example, we show the original image forward pass on the top and the augmented image forward pass on the bottom. From left to right: input images, image feature maps, un-augmented feature maps, and projected feature maps. Although the projected features are not the same between the original image and the augmented image, the feature distributions around the pedestrian locations (red points) remain coherent.

width. We thus include another bounding box regression loss \mathcal{L}_{box} using the same \mathcal{L}_1 distance as the offset loss \mathcal{L}_{off} over ground truth foot point locations.

Together, we have the overall loss functions for MVDetr,

$$\mathcal{L} = \mathcal{L}_{\text{det}} + \mathcal{L}_{\text{off}} + \frac{1}{C} \sum_c (\mathcal{L}_{\text{det},c} + \mathcal{L}_{\text{off},c} + 0.1 \times \mathcal{L}_{\text{box},c}), \quad (6)$$

where $\mathcal{L}_{\cdot,c}$ denotes the image-wise loss for a certain camera c .

View-Coherent Data Augmentation. Data augmentation helps to prevent the system from overfitting. However, this is non-trivial as multiview systems must maintain view consistency after augmentation. In this work, we design a view-coherent data augmentation scheme to meet such requirements.

First, for each input view, we apply random data augmentations including horizontal flipping, cropping, and scaling. Specifically, we record the affine matrix for the random transformations applied during data augmentation.

Next, we extract features from such augmented images. The per-view supervisions are also adjusted accordingly to the augmentation applied to the image for pedestrian detection.

Then, we apply the inverse affine transformation to create un-augmented feature maps (see Fig. 3). Such un-augmented feature maps do not equal the feature maps of the original image forward pass, but maintain the pedestrian locations and have similar feature distributions around these locations.

Finally, we project the un-augmented image features to the ground plane, using the same perspective transformations with

the camera pose and calibration. In this manner, we apply random data augmentation to the multiview images without breaking the view-consistency requirements.

4 EXPERIMENTS

4.1 Experiment Settings

Datasets. We verify the performance of the proposed detector and training scheme on the following datasets.

Wildtrack [4] is a real-world dataset that focuses on a 12×36 square meters region with 7 synchronized cameras. The ground plane is discretized into and annotated at a resolution of 480×1440 , and the images are of 1080×1920 pixels. It consists of 400 frames, with the first 360 frames dedicated for training and the last 40 frames used for testing.

MultiviewX [14] is a synthetic dataset created using the Unity [34] engine. It records a 16×25 square meters city square with 6 cameras, while having almost double the pedestrian crowdedness compared to the Wildtrack dataset to compensate for the limited appearance variance in the synthetic engine. Its ground plane is discretized into a 640×1000 grid, and the images are also taken with a 1080×1920 resolution. Similarly, it also has 400 frames with the last 40 frames used for evaluation.

Evaluation metrics. Different from monocular-view detection systems that evaluate the estimated bounding boxes, multiview detection systems evaluate the estimated ground plane occupancy

Method	Wildtrack				MultiviewX			
	MODA	MODP	Precision	Recall	MODA	MODP	Precision	Recall
RCNN & clustering [40]	11.3	18.4	68	43	18.7	46.4	63.5	43.9
POM-CNN [10]	23.2	30.5	75	55	-	-	-	-
DeepMCD [3]	67.8	64.2	85	82	70.0	73.0	85.7	83.3
Deep-Occlusion [1]	74.1	53.8	95	80	75.2	54.7	97.8	80.2
MVDeTr [14]	88.2	75.7	94.7	93.6	83.9	79.6	96.8	86.7
MVDeTr (w/ augmentation)	89.0	75.5	93.5	95.5	85.6	78.0	96.4	89.2
MVDeTr (deform attention)	90.2	81.6	96.9	93.2	93.0	91.1	99.4	93.6
MVDeTr (w/o augmentation)	89.5	81.2	94.0	95.6	93.1	90.9	98.5	94.4
MVDeTr (w/o per-view loss)	89.9	81.8	97.7	92.1	92.6	90.6	99.3	93.3
MVDeTr (ours)	91.5	82.1	97.4	94.0	93.7	91.3	99.5	94.2

Table 1: Performance (%) comparison on Wildtrack and MultiviewX datasets. “MVDeTr (ours)” is the proposed method. “MVDeTr (convolution)”, “MVDeTr (deform attention)”, “MVDeTr (w/o augmentation)”, and “MVDeTr (w/o per-view loss)” are variants of the proposed method. See Section 4.3 for more details.

map. Specifically, rather than intersection-over-union (IOU) between the estimated and ground truth bounding boxes, multiview detection systems consider the distance between the estimated pedestrian location and its ground truth, and use a 0.5-meter threshold for deciding true positives. Following previous works [4, 14], we use multi-object detection accuracy (MODA) as the primary indicator, since it considers both false positives and false negatives.

Implementation details. We first downsample all input images to 720×1280 for all $C = 7$ or $C = 6$ cameras since all images need to be considered at the same time. Inside the ResNet-18 feature extractor, we reduce the feature map downsample rate to 8 by replacing the last 2-strided convolutions with dilated convolutions, so as to increase its size while maintaining the receptive field. In this manner, we can preserve a relatively high-definition feature map, which benefits the detection system. To further reduce the memory footprint, we add a 128 channel bottleneck to the ResNet-18 feature extractor, before the per-view output heads and the feature projection module. For feature map perspective transformation and the ground plane heatmap supervision, we downsample the world grid by a factor of $r = 4$. For the shadow transformer, we build the transformer encoder for multiview processing and fusing with 3 encoder layers, each has $M = 8$ heads and attends to $K = 4$ points in each head. The shadow transformer also takes an additional $2\times$ downsampling on the projected ground plane feature maps due to memory concerns.

During training, we randomly apply horizontal flipping, cropping, and scaling with the proposed view-coherent augmentation. We train the detector using an Adam optimizer [15] with one-cycle learning rate scheduler [28] at a maximum learning rate of 1×10^{-3} . We initialize the ResNet-18 using the ImageNet [6] pre-trained model and set its learning rate to $0.1\times$ of all other layers. All experiments are conducted using one RTX-2080TI GPU.

4.2 Evaluation of MVDeTr

In Table 1, we compare the multiview detection performance of the proposed method against previous methods. Specifically, we find that MVDeTr outperforms the previous state-of-the-art by +3.3% MODA on the Wildtrack dataset, and +9.8% MODA on the

MultiviewX dataset. We point out that such improvements bring the MODA performance of MVDeTr to relatively high accuracies of 91.5% and 93.7% MODA on the two datasets.

In Fig. 4, we visualize the learned multiview deformable attention. All attended points ($M \times K \times C$ reference points in Eq. 3) of the same color are from a certain point of a certain camera. When considering the points that contribute mostly (in each circle), we find that the attention displays distinct patterns at different positions and different camera, while considering both their neighbors and the projected human bodies. Such adaptive attention patterns well align with our goal to deal with the various shadow-like distortion patterns, verifying the effectiveness of the proposed method.

We further visualize the detection results from the proposed method in Fig. 5, where most pedestrians are well detected even though the scene is heavily crowded and most pedestrians are occluded in certain views.

4.3 Ablations

Improvements over convolution. In order to quantitatively compare the performance difference between convolution-based multiview aggregation and the proposed shadow transformer alternative, we evaluate the convolution variant “MVDeTr (convolution)” in Table 1. “MVDeTr (convolution)” is very similar to the previous method MVDeTr [14]. Their main differences are that the former has a few modifications (e.g., a 128-dim bottleneck, different output heads) to fit it into one RTX 2080TI GPU, and adopts the same training scheme as the full MVDeTr. Replacing multiview deformable attention with convolution leads to -2.0% and -1.0% MODA losses on Wildtrack and MultiviewX, respectively. The raw MODA differences are relatively small, but we point out that they are over some very high baselines (the convolution variant) and are statistically very significant (*i.e.*, p -value < 0.001 over 5 runs). Although this variant achieves higher precisions, its lower recalls prevent it from outperforming the proposed method. We believe the main reason for this is that the convolution approach applies the same calculation over entire projected feature maps, which can provide slightly more confident results (higher precision) but often at the price of failing to detect at some uncertain positions (much lower recall).

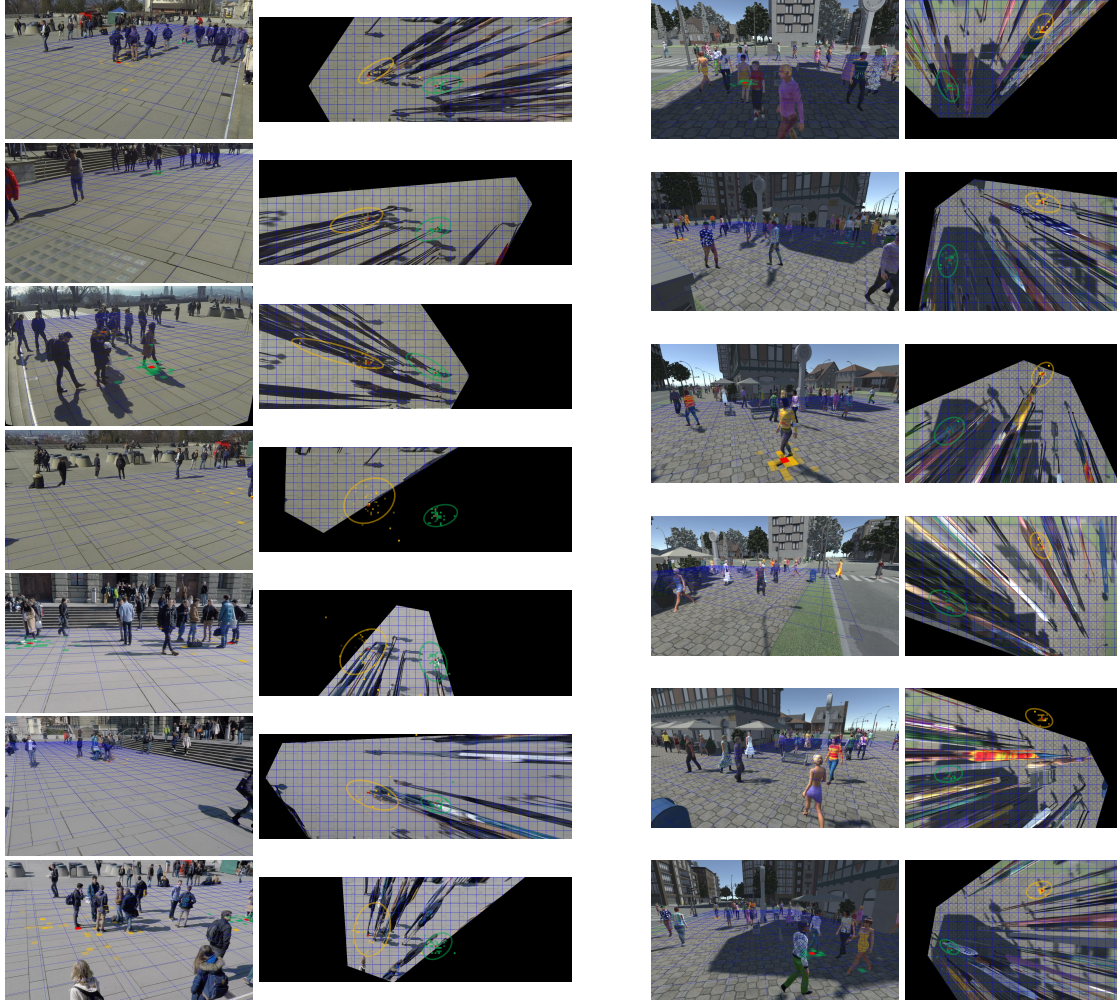


Figure 4: Visualization of the attended points and their corresponding weights in MVDeTr. We show the results on Wildtrack and MultiviewX datasets. For each dataset, we show the original image on the left, and the projected image onto the ground plane on the right. The learned attentions and attended positions on certain positions (red points) are shown with yellow or green points of different shades. We circle out the most prominent points (with high attention weights) in each projected view. The shadow transformer in MVDeTr learns to attend to different points at different positions and cameras, demonstrating its effectiveness in dealing with the various shadow-like distortion patterns.

Improvements over deformable attention. We show the effectiveness of our multiview deformable attention over deformable attention [47] in Table 1. Specifically, we replace all multiview deformable attention with deformable attention in our shadow transformer, and create the “MVDeTr (deform attention)” variant. In this manner, the multiview detector only fuses from the last fully connected layer before the output heads. When compared against the proposed method, “MVDeTr (deform attention)” falls short by -1.7% MODA on Wildtrack and -0.7% MODA on MultiviewX. Again, such performance differences are small in raw values but statistically significant, with p -value < 0.05 over 5 runs. This also verifies that our multiview deformable attention design is more effective for multiview detection, as it jointly considers multiple views to process and fuse.

Effectiveness of the view-coherent data augmentation. In Table 1, removing data augmentation from MVDeTr causes a -2.0% MODA loss on Wildtrack dataset and a -0.6% MODA loss on MultiviewX dataset. Adding the proposed view-coherent data augmentation also helps the MVDet system, and improves its MODA by $+0.8\%$ and $+1.7\%$ on two datasets. Overall, data augmentation gives better precision (fewer false positives) and worse recall (more false negatives). The system trained with data augmentation is less likely to consider an unsure position occupied by a pedestrian as it predicts fewer positives. Data augmentation might crop out a certain part of the image, which can create more unsure candidates. As the ground truth positives (on the ground plane) remain the same, more unsure candidates during training can reduce the tendency of considering unsure positions as occupied during testing, further



Figure 5: Detection result visualization of the proposed method.

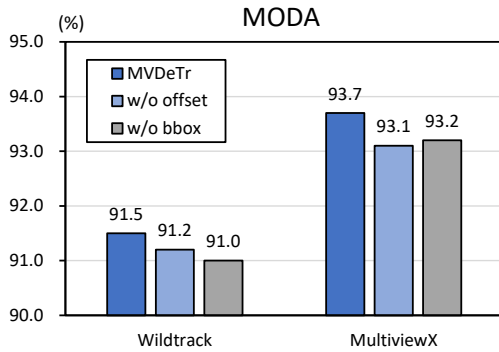


Figure 6: Influence of different loss terms.

leading to fewer predicted positives. On the whole, as the trade-off between fewer false positives and more false negatives still leads to higher MODA, we consider the view-coherent augmentation beneficial to the overall system.

Effectiveness of the per-view loss. As shown in Table 1, the removal of the per-view loss results in MODA drops of -1.6% and -1.1% on Wildtrack and MultiviewX, respectively. This aligns with previous findings by Hou *et al.* [14]. In fact, we believe that despite being formulated as a ground plane key-point detection problem, multiview detection can still benefit from the supervision of un-projected images. Especially, with the help of the view-coherent data augmentation, our correspondingly adjusted per-view supervision can also further prevent the system from overfitting.

Influence of the offset and bounding box loss terms. In Fig. 6, we further investigate the influence of the offset loss and the bounding box loss. Specifically, we find the removal of either of these leads to a performance drop on both datasets. Offset loss is slightly more important for the MultiviewX dataset (the removal of it causes a -0.6% MODA loss compared to -0.3% on Wildtrack), as the synthetic human can sometimes move too close together and the absence of the offset estimation can lead to potentially more false negatives from non-maximum suppression. Bounding

box regression, on the other hand, is similarly important for both datasets (the removal of it results in -0.5% MODA loss for both datasets), indicating that regressing the bounding box shapes is beneficial to the system training.

5 CONCLUSION

In this paper, we investigate how to jointly consider multiple views to combat the occlusions in detection. Specifically, we find that convolution might not be effective for multiview feature aggregation under various *position-sensitive* distortions due to its *translation-invariant* nature. In this case, we propose MVDeTr, a multiview detector that adopts shadow transformer for multiview feature fusion. Shadow transformer is designed specifically for the shadow-like distortions, and can attend accordingly at different positions over multiple camera views. Together with MVDeTr, we adopt a new training scheme that includes view-coherent data augmentation, which preserves multiview consistency after random augmentation. Using the novel training scheme and MVDeTr, we report new state-of-the-art on multiview detection benchmarks.

ACKNOWLEDGEMENT

This work was supported by the ARC Discovery Early Career Researcher Award (DE200101283) and the ARC Discovery Project (DP210102801).

REFERENCES

- [1] Pierre Baqué, François Fleuret, and Pascal Fua. 2017. Deep occlusion reasoning for multi-camera multi-target detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 271–279.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European Conference on Computer Vision*. Springer, 213–229.
- [3] Tatjana Chavdarova et al. 2017. Deep multi-camera people detection. In *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 848–853.
- [4] Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Timur Bagautdinov, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. 2018. Wildtrack: A multi-camera hd dataset for dense unscripted pedestrian detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5030–5039.

- [5] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. 2017. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*. 764–773.
- [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=YicbFdNTTy>
- [9] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. 2019. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 6569–6578.
- [10] Francois Fleuret, Jerome Berclaz, Richard Lengagne, and Pascal Fua. 2007. Multi-camera people tracking with a probabilistic occupancy map. *IEEE transactions on pattern analysis and machine intelligence* 30, 2 (2007), 267–282.
- [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 580–587.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [13] Jan Hosang, Rodrigo Benenson, and Bernt Schiele. 2017. Learning non-maximum suppression. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4507–4515.
- [14] Yunzhong Hou, Liang Zheng, and Stephen Gould. 2020. Multiview detection with feature perspective transformation. In *European Conference on Computer Vision*. Springer, 1–18.
- [15] Diederik P Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *ICLR (Poster)*.
- [16] Hei Law and Jia Deng. 2018. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 734–750.
- [17] Jinpeng Li, Shengcai Liao, Hangzhi Jiang, and Ling Shao. 2020. Box Guided Convolution for Pedestrian Detection. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 1615–1624. <https://doi.org/10.1145/3394171.3413989>
- [18] Wei Li, Zhenting Wang, Xiao Wu, Ji Zhang, Qiang Peng, and Hongliang Li. 2020. CODAN: Counting-Driven Attention Network for Vehicle Detection in Congested Scenes. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 73–82. <https://doi.org/10.1145/3394171.3413945>
- [19] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*. 2980–2988.
- [20] Wei Liu, Shengcai Liao, Weiqiang Ren, Weidong Hu, and Yinan Yu. 2019. High-level semantic feature detection: A new perspective for pedestrian detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5187–5196.
- [21] Junhyug Noh, Soochan Lee, Beomsu Kim, and Gunhee Kim. 2018. Improving occlusion and hard negative handling for single-stage pedestrian detectors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 966–974.
- [22] Wanli Ouyang, Xingyu Zeng, and Xiaogang Wang. 2015. Partial occlusion handling in pedestrian detection with a deep model. *IEEE Transactions on Circuits and Systems for Video Technology* 26, 11 (2015), 2123–2137.
- [23] Zengyi Qin, Jinglu Wang, and Yan Lu. 2020. Weakly Supervised 3D Object Detection from Point Clouds. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 4144–4152. <https://doi.org/10.1145/3394171.3413805>
- [24] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*. 91–99.
- [25] Hamid Rezaatfighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. 2019. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 658–666.
- [26] Gemma Roig, Xavier Boix, Horesh Ben Shitrit, and Pascal Fua. 2011. Conditional random fields for multi-camera object detection. In *2011 International Conference on Computer Vision*. IEEE, 563–570.
- [27] Bing Shuai, Andrew Berneshawi, Manchen Wang, Chunhui Liu, Davide Modolo, Xinyu Li, and Joseph Tighe. 2020. Application of Multi-Object Tracking with Siamese Track-RCNN to the Human in Events Dataset. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 4625–4629. <https://doi.org/10.1145/3394171.3416297>
- [28] Leslie N Smith and Nicholay Topin. 2019. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, Vol. 11006. International Society for Optics and Photonics, 1100612.
- [29] Tao Song, Leiyu Sun, Di Xie, Haiming Sun, and Shiliang Pu. 2018. Small-scale pedestrian detection based on topological line localization and temporal feature aggregation. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 536–551.
- [30] Maarten Sukel, Stevan Rudinac, and Marcel Worring. 2020. Detecting Urban Issues With the Object Detection Kit. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 4518–4520. <https://doi.org/10.1145/3394171.3414427>
- [31] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. 2019. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7464–7473.
- [32] Yonglong Tian, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep learning strong parts for pedestrian detection. In *Proceedings of the IEEE international conference on computer vision*. 1904–1912.
- [33] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. 2019. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE international conference on computer vision*. 9627–9636.
- [34] Unity. [n.d.]. Unity Technologies. <https://unity.com/>
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>
- [36] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. 2018. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7794–7803.
- [37] Xinlong Wang, Tete Xiao, Yuning Jiang, Shuai Shao, Jian Sun, and Chunhua Shen. 2018. Repulsion loss: Detecting pedestrians in a crowd. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7774–7783.
- [38] Jialian Wu, Chunluan Zhou, Qian Zhang, Ming Yang, and Junsong Yuan. 2020. Self-Mimic Learning for Small-Scale Pedestrian Detection. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 2012–2020. <https://doi.org/10.1145/3394171.3413634>
- [39] Xiongwei Wu, Doyen Sahoo, and Steven Hoi. 2020. Meta-RCNN: Meta Learning for Few-Shot Object Detection. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 1679–1687. <https://doi.org/10.1145/3394171.3413832>
- [40] Yuanlu Xu, Xiaobai Liu, Yang Liu, and Song-Chun Zhu. 2016. Multi-view people tracking via hierarchical trajectory composition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4256–4265.
- [41] Fuzhi Yang, Huan Yang, Jianlong Fu, Hongtao Lu, and Baining Guo. 2020. Learning texture transformer network for image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5791–5800.
- [42] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. 2019. Cross-modal self-attention network for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10502–10511.
- [43] Cai YuanQiang, Dawei Du, Libo Zhang, Longyin Wen, Weiqiang Wang, Yanjun Wu, and Siwei Lyu. 2020. Guided Attention Network for Object Detection and Counting on Drones. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 709–717. <https://doi.org/10.1145/3394171.3413816>
- [44] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z Li. 2018. Occlusion-aware R-CNN: detecting pedestrians in a crowd. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 637–653.
- [45] Jiewen Zhao, Ruize Han, Yiyang Gan, Liang Wan, Wei Feng, and Song Wang. 2020. Human Identification and Interaction Detection in Cross-View Multi-Person Videos with Wearable Cameras. In *Proceedings of the 28th ACM International Conference on Multimedia (Seattle, WA, USA) (MM '20)*. Association for Computing Machinery, New York, NY, USA, 2608–2616. <https://doi.org/10.1145/3394171.3413903>
- [46] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. 2019. Objects as points. *arXiv preprint arXiv:1904.07850* (2019).
- [47] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. 2021. Deformable (DETR): Deformable Transformers for End-to-End Object Detection. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=gZ9hCDWe6ke>