

Booster-SHOT: Boosting Stacked Homography Transformations for Multiview Pedestrian Detection with Attention

Jinwoo Hwang*
 Seoul National University
 luorix@snu.ac.kr

Philipp Benz
 Deeping Source Inc.
 philipp.benz@deepingsource.io

Pete Kim
 Deeping Source Inc.
 pete.kim@deepingsource.io

Abstract

*Improving multi-view aggregation is integral for multi-view pedestrian detection, which aims to obtain a bird’s-eye-view pedestrian occupancy map from images captured through a set of calibrated cameras. Inspired by the success of attention modules for deep neural networks, we first propose a Homography Attention Module (HAM) which is shown to boost the performance of existing end-to-end multiview detection approaches by utilizing a novel channel gate and spatial gate. Additionally, we propose Booster-SHOT, an end-to-end convolutional approach to multiview pedestrian detection incorporating our proposed HAM as well as elements from previous approaches such as view-coherent augmentation or stacked homography transformations. Booster-SHOT achieves 92.9% and 94.2% for MODA on Wildtrack and MultiviewX respectively, outperforming the state-of-the-art by 1.4% on Wildtrack and 0.5% on MultiviewX, achieving state-of-the-art performance overall for standard evaluation metrics used in multi-view pedestrian detection.*¹

1. Introduction

Multi-view detection [2, 26, 41] leverages multiple camera views for object detection using synchronized input images captured from varying view angles. Compared to a single-camera setup, the multi-view setup alleviates the occlusion issue, one of the fundamental problems in many computer vision applications. In this work, we consider the problem of multi-view pedestrian detection. As shown in Figure 1, a bird’s-eye-view representation is obtained with the synchronized images from multiple calibrated cameras, which is then further used to detect pedestrians in the scene.

A central problem in multi-view detection is to obtain a correct *multi-view aggregation*. The change in viewpoint and occlusions make it challenging to match object features

across different view angles. Various works attempted to address this problem, ranging from early approaches leveraging “classical” computer vision [3], hybrid approaches further incorporating deep learning, to end-to-end trainable deep learning architectures [25, 26, 42].

One core challenge in multiview detection is designing how the multiple views should be aggregated. MVDet [26] proposes a fully convolutional end-to-end trainable solution for the multi-view detection task. MVDet aggregates different views by projecting the convolution feature map via perspective transformation to a single ground plane and concatenating the multiple projected feature maps. Given the aggregated representation, MVDet applies convolutional layers to detect pedestrians in the scene. Song *et al.* [42] identified that the projection of the different camera views to a single ground plane is not accurate due to misalignments. Consequently, they proposed to project the feature maps onto different height levels according to different semantic parts of pedestrians. Additionally, they use a neural-network-based soft-selection module to assign a likelihood to each pixel of the features extracted from the different views. They termed their approach SHOT, due to the use of the Stacked HOMography Transformations. MVDeTr [25] extends MVDet by introducing a shadow transformer to attend differently at different positions to deal with various shadow-like distortions as well as a view-coherent data augmentation method, which applies random augmentations while maintaining multiview-consistency. MVDeTr currently constitutes the SotA approach for multiview detection.

In recent years the attention mechanism for deep neural networks has played a crucial role in deep learning [21, 27, 46] due to the non-trivial performance gains that it enabled. Attention mechanisms have provided benefits for various vision tasks, e.g. image classification [27, 46], object detection [7, 13], semantic segmentation [17, 51], or Point Cloud Processing [20, 47]. However, to this date, no dedicated attention mechanism has been proposed for the task of multiview pedestrian detection.

In this work, we fill this gap and propose an attention

*This work was performed while Jinwoo was at Deeping Source Inc.

¹Code can be found at <https://github.com/luorix1/Booster-SHOT>.

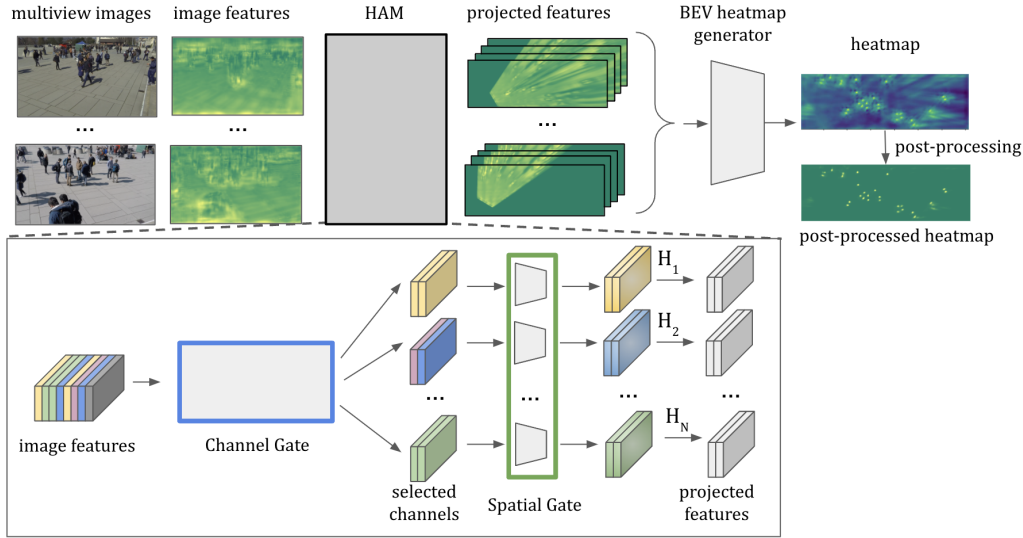


Figure 1. Overview of multiview detection with homography attention module (HAM)

mechanism specifically designed to boost existing multiview detection frameworks. Our proposed Homography Attention Module (HAM) is specifically tailored for the core task of multiview aggregation in modern multiview detection frameworks. As shown in the lower part of Figure 1 our proposed solution consists of a channel gate module and a spatial gate module. The channel gate is directly applied to the accumulated image features from the different views. The intuition behind our *channel gate* is that different channels hold meaningful information for different homographies. The channel gate is followed by our *spatial gate*. We conjecture, that for each view and homography combination different spatial features are of higher importance. Our proposed attention mechanism can be readily plugged into existing methods.

We also combine insight from previous approaches and HAM to propose Booster-SHOT, a new end-to-end multiview pedestrian detection framework. Our experimental results show that both incorporating HAM into previous frameworks and Booster-SHOT improves over previous multiview detection frameworks and achieves state-of-the-art performance. Additionally, we provide quantitative and qualitative results to verify and justify our design choices.

2. Related work

2.1. Multiview Detection

Detecting pedestrians in crowded and occluded scenes is challenging with a single camera view. Hence, researchers focus on multi-camera setups, which provide a richer representation of the environment. Calibrated and synchronized cameras establish correspondences between ground

plane locations and bounding boxes in multiple camera views. Early methods for multiview pedestrian detection used background subtraction, geometric constraints, and occlusion reasoning [6, 16, 41]. Fleuret *et al.* [16] estimated a probabilistic occupancy map and performed tracking, while Sankaranarayanan *et al.* [41] leveraged geometric constraints. Coates and Ng [12] fused object detector outputs probabilistically. Berclaz *et al.* [6] proposed a multiple object tracking framework using the k-shortest paths algorithm, and Roig *et al.* [39] modeled the problem using Conditional Random Fields.

Deep neural networks have been successfully applied to multi-view pedestrian detection [5, 9, 25, 26, 42]. Chavdarova and Fleuret [9] proposed an end-to-end architecture combining monocular pedestrian detectors with a multi-view network. Baqu’*et al.* [5] addressed performance degradation in crowded scenes using a hybrid CRF-CNN approach.

MVDeT [26] is an end-to-end trainable multiview detector that aggregates cues by transforming feature maps to a single ground plane and using large kernel convolutions. Extensions to MVDeT include stacked homographies and shadow transformers [25, 42]. Song *et al.* [42] improved alignment using 3D world coordinate projections and introduced a soft selection module. MVDeTr [25] adopted shadow transformers, which attend differently based on position and camera differences, and introduced view-coherent data augmentation.

The recent works of Lee *et al.* [31], Engliberge *et al.* [15], and Gao *et al.* [18] have further advanced multiview pedestrian detection. Lee *et al.* proposed Multiview Target Transformation (MVTT) to address distortion caused

by perspective transformation. Engliberge *et al.* [15] introduced a novel multiview data augmentation pipeline to preserve alignment among views. Gao *et al.* [18] incorporated key point supervision and grouped feature fusion to enhance multiview pedestrian detection.

2.2. Attention Mechanism in Computer Vision

Attention mechanisms for computer vision emphasize more important regions of an image or a feature map and suppress less relevant parts [21]. They can be broadly divided into channel attention, spatial attention, and a combination of the two variants.

Channel Attention selects important channels through an attention mask across the channel domain. Pioneered by Hu *et al.* [27] various works have extended upon the Squeeze-and-Excitation (SE) mechanism module [19, 30, 38, 49].

Spatial Attention selects important spatial regions of an image or a feature map. Early spatial attention variants are based on recurrent neural networks (RNN) [4, 35]. In the literature various variants of visual attention-based model can be found [36, 48] To achieve transformation invariance while letting CNNs focus on important regions, Spatial Transformer Networks [28] had been introduced. Similar mechanisms have been introduced in deformable convolutions [13, 52]. Originating from the field of natural language processing, self-attention mechanisms have been examined for computer vision applications [7, 11, 14, 45, 53].

Channel Attention & Spatial Attention can also be used in combination. Residual Attention Networks [44] extend ResNet [22] through a channel & spatial attention mechanism on the feature representations. A spatial and channel-wise attention mechanism for image captioning has been introduced in [10]. The Bottleneck Attention Module (BAM) [37] and Convolutional Block Attention Module (CBAM) [46] both infer attention maps along the channel and spatial pathway. While in the previous two methods the channel and spatial pathways are computed separately, triplet attention [34] was introduced to account for cross-dimension interaction between the spatial dimensions and channel dimension of the input. Channel & spatial attention has also been applied in the context of segmentation [17, 40] Further combinations of channel and spatial attention include self-calibrated convolutions [32], coordinate attention [24] and strip pooling [23].

3. Preliminaries

Let the input images for N camera views be (I^1, \dots, I^N) . The respective feature maps obtained from the feature extractor in the initial step of the general framework are denoted as (F^1, \dots, F^N) . The intrinsic, extrinsic parameters of the i 'th camera are $\mathbf{G}^i \in \mathbb{R}^{3 \times 3}$ and $\mathbf{E}^i = [\mathbf{R}^i | \mathbf{t}^i] \in \mathbb{R}^{3 \times 4}$, respectively, where \mathbf{R}^i is the 3×3 matrix for rota-

tion in the 3D space and \mathbf{t}^i is the 3×1 vector representing translation. Following MVDet [26], we quantize the ground plane into grids and define an additional matrix $\mathbf{F}^i \in \mathbb{R}^{3 \times 3}$ that maps world coordinates to the aforementioned grid. While the mathematical concept of homography is an isomorphism of projective spaces, we use the term homography to describe correspondence relations between points on a given plane parallel to the ground as they are seen from the bird's-eye-view and from a separate camera-view. This is in line with SHOT [42] where the authors explain their projections as being homographies describing the translation of a plane for the pin-hole camera model. We will go into further depth regarding the homography transforms in our supplementary materials.

4. Methodology

4.1. Previous Multiview Detection Methods

Before presenting our proposed attention module, we outline the previous multiview detection frameworks on which this work builds upon. MVDet [26] presented a multiview detection framework that functions as follows: First, the input images from different viewpoints are passed through a generic feature extractor such as ResNet18 with minor modifications. The feature maps are passed through an additional convolutional neural network that detects the head and feet of pedestrians to aid the network during training. Next, the feature maps are projected to the ground plane via homography transformation and concatenated. Additionally, x, y coordinate maps are concatenated to the stack of transformed feature maps as in CoordConv [33]. Finally, this is passed through a CNN to output a bird's-eye-view (BEV) heatmap which is then post-processed via thresholding and non-maximum suppression. Extending upon MVDet, MVDeTr [25] proposed the use of affine transformations, which are view-coherent augmentations. Additionally, the final CNN to generate the BEV heatmap is replaced with a shadow transformer, with the purpose to handle various distortion patterns during multiview aggregation. MVDeTr further replaces the MSE loss used in MVDet with Focal Loss [29] coupled with an offset regression loss. While MVDet and MVDeTr both project the feature maps to the ground plane, SHOT [42] proposes to approximate projections in 3D world coordinates via a stack of homographies. In line with MVDet, SHOT uses a ResNet18 as a feature extractor. Contrary to MVDet, SHOT introduces additional planes parallel to the ground plane with different distances to the ground. The features are selectively projected from the camera-view to these different bird's-eye-views. As the projection to some planes may be of more importance than others, SHOT introduces a soft selection module where a network learns which homography should be used for which pixel.

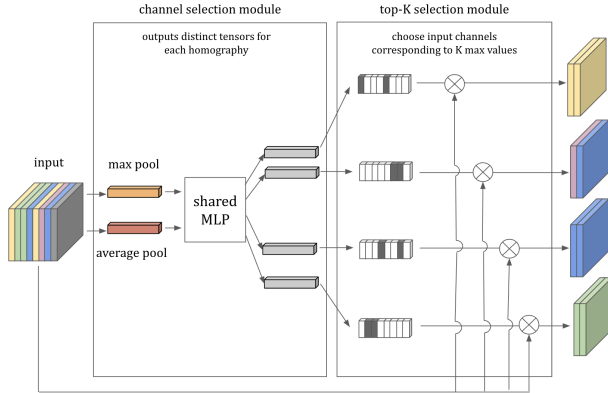


Figure 2. Diagram showing our proposed channel gate.

4.2. Homography Attention Module (HAM)

We identify two shortcomings of the soft selection module of SHOT [42]. First, it uses softmax activation, therefore each pixel gets projected to some extent to each homography. Since even the homography given the lowest score by the soft selection module affects the projected outcome, this introduces some noise into the final projected feature map. In addition, all feature channels corresponding to a single pixel are multiplied by the same value when projected to a homography. However, different channels attend to different features and some will be useful for the selected homography while others won't. For this reason we designed HAM such that it selects channels in a discrete manner for each homography to avoid the two problems mentioned above.

Here we outline our design choices for our proposed homography attention module (HAM) and their respective motivations. HAM consists of a channel gate and several spatial gates equal to the number of homographies used. Note, that our attention module is specifically designed for view-aggregation in the context of multiview detection and is hence only applied in the multiview aggregation part. The image feature maps are first passed through the channel gate, then the spatial gate, and then finally through the homography, followed by the BEV heatmap generator.

Channel Gate Our proposed channel gate follows the intuition that depending on the homography, different channels are of importance. Taking into consideration the multiple homography layers deployed at different heights, different feature information become more valuable. For instance, when we consider the homography at $Z = 0$, discriminative feature information near the ground plane, such as a person's feet, ankles, and lower legs, may offer more significant representation. This is because the homography at $Z = 0$ focuses on objects that are closer to the ground

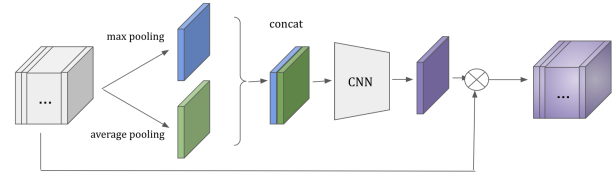


Figure 3. Diagram showing our proposed spatial gate.

plane, which makes features near the ground plane more informative. This is in contrast to the approach proposed by SHOT, which feeds all feature maps through each of the homographies. Figure 2 outlines the architecture of our proposed channel gate, which broadly consists of the *channel selection module* and the *top-K selection module*. Given the stack of feature maps acquired from the different views, first the channel selection module is applied. The channel selection module first applies max pooling and average pooling along the spatial dimension. Both pooled feature maps are passed through a shared 2-layer MLP. As the number of channels in the output from the last layer of the MLP is decided by the number of homographies (denoted as D) with the number of channels in the input (denoted as C), we obtain C channel for each homography, or in other words a $C \times D$ channel output size. Afterward, we apply the softmax function along the channel dimension for each of the outputs. The outputs are then fed into the top-K selection module. The top-K selection module takes these D different C -dimensional outputs and selects the top K largest values. The corresponding top-K selected channels from the original input are then concatenated, resulting in a subset of the original input with K channels. In the case of $D = 1$ (using only one homography, usually the ground plane), the top-K selection module defaults to an identity function. To retain the channel-wise aspect of our module, in this scenario we multiply the output of the channel selection module element-wise with the input. This completes the channel gate, which outputs D times K -channel feature maps, which are then fed into spatial gate.

Spatial Gate Our spatial gate is motivated by our conjecture that for each view and homography combination different spatial features are of different importance. This intuition is based on the understanding that the view determines the camera's position and orientation, the homography corresponds to different heights in the scene, and the spatial features capture the patterns, textures, and shapes of the objects in the scene. Depending on the specific view and homography combination, certain spatial features may be more informative and relevant for feature extraction than others. For example, features closer to the lower image border might be more important for a view-homography com-

Table 1. Settings for each approach

Method	Aug.	Loss	BEV gen.	Multi Homogr.
MVDeTr [26]	✗	MSE	CNN	✗
SHOT [42]	✗	MSE	CNN	✓
MVDeTr [25]	✓	Focal	Transformer	✗
Booster-Shot	✓	Focal	CNN, Transformer	✓

bination with a nearly parallel to the ground plane and the homography at $Z = 0$. By using a spatial gate to selectively weight and filter the spatial features for each combination, our proposed method can effectively capture the relevant information from the image and improve performance. Figure 3 shows the architecture of our spatial gate. The input is max and average pooled along the channel dimension, then concatenated channel-wise. This 2-channel input is then passed through a 2-layer convolutional neural network to generate the spatial attention map. Finally, this spatial attention map is multiplied with the original input element-wise to create an output with dimensions identical to the input. For each homography-path a separate spatial gate is applied.

Architecture-wise, while SHOT uses a “soft selection module” to estimate the importance of each homography plane for each pixel, HAM estimates the importance of channels and spatial information for each homography. Also, while MVDeTr introduced a “shadow transformer” after the homography transforms to remove shadow-like background noise, HAM uses attention to optimize image features fed to each homography and is applied prior to the homography transforms.

4.3. Booster-SHOT

Given the insights collected from previous approaches in addition to our proposed HAM, we design a multiview pedestrian detection architecture, which we term Booster-SHOT. Booster-SHOT bases itself on SHOT [42], using their stacked homography approach and leverages MVDeTr’s Focal loss and offset regression loss along with the view-coherent augmentation. We retain SHOT’s convolutional architecture used to generate the BEV heatmap but remove the soft selection module as the implementation of our module renders it obsolete. Figure 1 outlines how our proposed module is implemented in Booster-Shot. Table 1 outlines the design choices of Booster-SHOT alongside previous methods.

5. Experiments

5.1. Datasets

Our method is tested on two datasets for multiview pedestrian detection.

Wildtrack [8] consists of 400 synchronized image pairs from 7 cameras, constituting a total of 2,800 images. The images cover a region with dimensions 12 meters by 36 meters. The ground plane is denoted using a grid of dimensions 480×1440 , such that each grid cell is a 2.5-centimeter by 2.5-centimeter square. Annotations are provided at 2fps and there are, on average, 20 people per frame. Each location within the scene is covered by an average of 3.74 cameras.

MultiviewX [26] is a synthetic dataset created using human models from PersonX [43] and the Unity engine. It consists of 1080×1920 images taken from 6 cameras that cover a 16-meter by 25-meter area. Per the method adopted in Wildtrack, the ground plane is represented as a 640×1000 grid of 2.5-centimeter squares. Annotations are provided for 400 frames at 2fps. An average of 4.41 cameras cover each location, while an average of 40 people are present in a single frame.

5.2. Settings and metrics

In accordance with the previous methods, we report the four metrics: Multiple Object Detection Accuracy (MODA), Multiple Object Detection Precision (MODP), precision, and recall. Let us define N as the number of ground truth pedestrians. If the true positives (TP), false positives (FP) and false negatives (FN) are known, precision and recall can be calculated as $\frac{TP}{FP+TP}$ and $\frac{TP}{N}$, respectively. MODA is an accuracy metric for object detection tasks and is therefore obtained by calculating $1 - \frac{FP+FN}{N}$. MODP is computed with the formula $\frac{\sum 1-d[d<t]/t}{TP}$ where d is the distance from a detection to its ground truth (GT) and t is the threshold for a correct detection. We keep the original threshold of 20 that was proposed in SHOT. Our implementation is based on the released code for MVDeTr [26], SHOT [42], MVDeTr [25] and follows the training settings (optimizer, learning rate, etc.) for each. For all instances, the input images are resized to 720×1280 images. The output features are 270×480 images for MVDeTr and SHOT and 90×160 for MVDeTr. Δz (the distance between homographies) is set to 10 on Wildtrack and 0.1 on MultiviewX. All experiments are run on two A30 GPUs (depending on the framework) with a batch size of 1.

For experiments implementing our module in SHOT, our base approach involves selecting the top-32 channels each for 4 homographies. We note that SHOT’s base approach uses 5 homographies.

5.3. Comparison with previous methods

First, to show the efficiency of our homography attention module (HAM), we use the previous approaches as is without modification to their loss or training configuration, and simply plug in our proposed HAM. As the soft selection module in SHOT is rendered obsolete by our proposed HAM, we remove it when comparing the performance of

Table 2. Performance comparison (in %) on Wildtrack and MultiviewX datasets

Method	Wildtrack				MultiviewX			
	MODA	MODP	precision	recall	MODA	MODP	precision	recall
MVDeTr	88.2	75.7	94.7	93.6	83.9	79.6	96.8	86.7
MVDeTr + HAM	89.6 ± 0.35	80.4 ± 0.21	95.7 ± 1.06	93.8 ± 0.42	91.3 ± 0.35	81.7 ± 0.14	98.3 ± 0.49	91.9 ± 2.26
SHOT	90.2	76.5	96.1	94.0	88.3	82.0	96.6	91.5
SHOT + HAM	90.2 ± 0.49	77.4 ± 0.57	96.2 ± 0.07	93.9 ± 0.42	91.2 ± 0.53	86.9 ± 4.14	98.2 ± 1.25	92.9 ± 0.78
MVDeTr	91.5	82.1	97.4	94.0	93.7	91.3	99.5	94.2
MVDeTr + HAM	92.8 ± 0.49	82.4 ± 0.71	96.6 ± 0.85	96.6 ± 1.25	94.2 ± 0.07	91.4 ± 0.57	99.4 ± 0.21	94.8 ± 0.21
Booster-Shot + Tr	92.5 ± 0.64	82.0 ± 0.71	96.3 ± 0.78	96.3 ± 1.50	93.8 ± 0.49	91.8 ± 0.07	98.8 ± 0.64	95.0 ± 1.06
Booster-Shot	92.8 ± 0.17	84.9 ± 4.42	97.5 ± 1.25	95.3 ± 1.17	94.4 ± 0.18	92.0 ± 0.04	99.4 ± 0.07	94.9 ± 0.21

SHOT with our module with the reported values for SHOT.

As shown in Table 2, applying our module to MVDeTr, SHOT and MVDeTr improved (or matched) all four metrics reported in their respective papers for MultiviewX. Specifically, the average performance of MVDeTr with our module improves over the reported values for MVDeTr on MultiviewX by 7.4%, 2.1%, 1.5%, and 5.2% for MODA, MODP, precision, and recall respectively. For Wildtrack, the use of our module again improved all four metrics with the exception of MVDeTr. For MVDeTr, our precision was still comparable with the reported value as there was only a 0.8% decrease in precision while the MODA, MODP, and recall each improved 1.3%, 0.3%, and 2.6% respectively.

The proposed Booster-SHOT outperforms previous methods across all metrics except for precision against MVDeTr. As MVDeTr proposed the shadow transformer as a way to improve performance, we applied it to Booster-SHOT and the results are denoted in Table 2 as Booster-SHOT + Tr. However, we were unable to obtain any meaningful improvement over the purely convolutional approach.

5.4. Analysis

Efficacy of HAM in comparison to existing methods

We emphasize that the novelty of HAM lies in the architectural integration of the attention mechanism for the specific purpose of multi-view aggregation, for which, to the best of our knowledge, our work is the first. Previous attention mechanisms (e.g. CBAM [46], CCG [1]) are applied at the convolutional blocks in the backbone network, while HAM is applied after the backbone network since it is tailored toward multi-view aggregation. Consequently, HAM can be seen as complementary to existing attention mechanisms.

To illustrate the importance of the design choices of HAM we compare it with the naive integration of SENet, CBAM, and CCG into Booster-SHOT on MultiviewX. SENet, CBAM, and CCG come after the feature extractor in place of HAM. To provide a common baseline for HAM, SENet, CBAM, and CCG, we provide additional results for “BoosterSHOT without attention”. This implementation is

equivalent to SHOT [42] with Focal Loss and training-time augmentations.

As shown in Table 3, BoosterSHOT outperforms all of the compared methods across the board. Only Booster-SHOT without attention shows similar results in precision, a very saturated metric for which BoosterSHOT shows only a slightly lower performance. In addition, when compared with BoosterSHOT without attention, adding CBAM, CCG, and SE showed only an increase of a maximum 0.5% in MODA, while adding HAM boosted MODA by 1.2%.

Attention for different homographies We previously conjectured that the significance of each channel is different for each homography. In the following we validate this hypothesis through empirical evidence. Note that the following results are shown for the synthetic MultiviewX dataset. Figure 4 shows images created from camera view 1 of the MultiviewX dataset and containing output from the channel selection module corresponding to each homography. The channel selection module output is average pooled channel-wise (in this instance, the output for each homography contains 32 channels) and superimposed onto a grayscale version of the original image from the MultiviewX dataset. Yellow areas indicate high values in the output, indicating that the network is attending strongly to those regions. We denote the ground plane as H0 (homography 0) and number the remaining homographies accordingly.

We can observe that the output from the channel selection module is homography-dependent as the yellow areas in all four images differ. We also note that the body parts with the brightest colors align with the height of the homographies. H0 highlights the feet while H1 highlights the lower body, especially around the knee area. H2 and H3 both highlight the upper body but H3 extends a bit farther upwards compared to H2. A similar phenomenon has been reported by the SHOT authors for their soft selection module. However, our channel selection module output shows more distinct highlighting of the body parts. Overall, these results support the importance of selecting different chan-

Table 3. BoosterSHOT performance with HAM vs pre-existing attention mechanisms

	MODA	MODP	precision	recall
Booster-SHOT w/o attention	93.2 ± 0.18	91.2 ± 0.07	99.4 ± 0.04	93.7 ± 0.20
Booster-SHOT (SE)	93.7 ± 0.23	88.2 ± 5.66	98.1 ± 1.11	95.5 ± 0.90
Booster-SHOT (CBAM)	93.2 ± 0.14	90.5 ± 0.14	98.5 ± 0.53	94.7 ± 0.35
Booster-SHOT (CCG)	93.4 ± 0.18	91.4 ± 0.04	99.1 ± 0.11	94.2 ± 0.07
Booster-SHOT	94.4 ± 0.18	92.0 ± 0.04	99.4 ± 0.07	94.9 ± 0.21

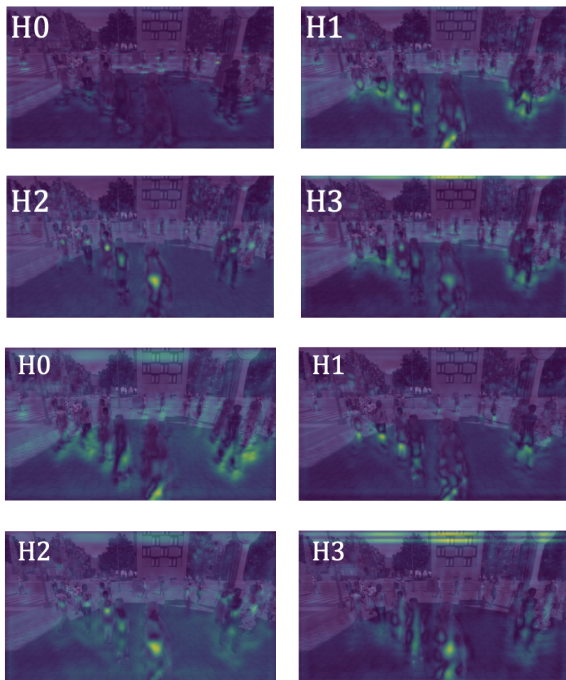


Figure 4. Homography-wise output from channel selection (top) and spatial attention maps (bottom)

nels for different homographies.

Figure 4 (bottom) shows the attention values from the spatial attention block at the end of our proposed module. All four attention maps show starkly different distributions, confirming our conjecture that different pixels in the feature map can differ in importance for each homography.

The results demonstrated above were obtained through an experiment where the distance between homography planes was increased from 10cm to 60cm for MultiviewX. Due to the low height of even the top homography plane in the 10cm case (30cm off the ground), the difference between the attention module outputs was not easily noticeable. By increasing the distance between homography planes, we were able to obtain images that clearly show homographies that are higher off the ground attend to higher regions of the human body. In addition, we noticed that the foot regression auxiliary loss caused bias toward the foot

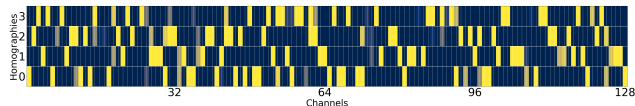


Figure 5. Heatmap representation of channel selection homography-wise. Deeper yellow colors indicates that the channel is selected most of the time while deeper blue colors are assigned to channels that are seldom selected.

region in the extracted image features, thus distorting our heatmap visualization of the attention module outputs. As such, the experiments from which Figure 4, Figure 5 and Figure 6 were obtained did not include auxiliary losses during training. For a more detailed analysis regarding auxiliary losses and their effect on performance, we refer the reader to our supplementary.

We further provide results averaged over the entire MultiviewX test dataset. Specifically, we visualize how often certain channels are selected for each homography for a given view for Booster-SHOT. For each channel, we count the number of times it is selected for each homography and divide by the total number of image pairs in the test set and display the resulting heatmap in Figure 5.

First, it can be observed that the channels that are selected often (yellow hues) show almost no overlap across homographies, again providing evidence to our previous claim that different channels attend to different homographies. Although there are minor differences in the specific number of times some channels are chosen, the channels that are selected for the majority of the test set for each homography are unchanged. (Due to length constraints, we refer the reader to our supplementary.) Interestingly, we also observe that some channels are not selected at all by any homography while other channels appear to be selected by multiple homographies.

Attention across different views Figure 6 further presents evidence that our channel selection module output is only homography-dependent. We denote the camera views as C1 (camera 1) through C6 for the homography to the ground plane (H0). For all 6 images, the feet and surrounding area of the pedestrians are highlighted. Therefore,



Figure 6. Camera view-wise output from channel selection module

we conclude that the output from the channel selection module attends consistently to certain features across all camera views.

Generalization across camera views. Due to the specific nature of camera pose and environmental factors, multiview pedestrian detection methods are susceptible to overfitting to a single scene. To test generalization capabilities, we adapt Table 4 in SHOT [42] and train models on a split of cameras of MultiviewX and evaluate the remaining cameras. Specifically, the top-view grid of the ground plane is divided into two separate grids of equal area. As MultiviewX has 6 cameras, 3 cameras cover each of the newly formed grids. We compare our approach to SHOT and MVDet.

The results presented in Table 4 show that introducing HAM results in significant improvement in MODA and precision for MVDet and SHOT, with MVDet’s MODA increasing 27.4% and SHOT’s MODA increasing 14.7%. Although a similar increase in MODP and recall was not observed, both metrics remain comparable to MVDet and SHOT. Through these results, HAM is shown to boost the generalization ability of both a single-homography (MVDet) approach and a multi-homography (SHOT) approach.

Table 4. Results from train-test camera split scenario on MultiviewX

	MODA	MODP	precision	recall
MVDet	33.0	76.5	64.5	73.4
MVDet + HAM	60.4	75.2	85.6	72.7
SHOT	49.1	77.0	73.3	77.1
SHOT + HAM.	63.8	76.6	86.0	76.2

Backbone network As per the methodology proposed in MVDet [26], we deploy ResNet-18 [22] as the backbone architecture. To keep the feature map resolution at a sufficiently high resolution, we substitute the final 2 strided convolutions with dilated convolutions [50]. To minimize memory usage, a bottleneck is added to the backbone architecture with 128 channels. To ablate the influence of the backbone architecture, we compare this setting with an SE-

Net backbone, which we initialize with the ImageNet pre-trained versions from the TIMM library.² For a fair comparison, we used BoosterSHOT with the shadow transformer proposed in MVDeTr so the only difference is the use of HAM with multiple homographies. In this case, we used 4 homographies for BoosterSHOT (Tr) noted in Table 5. Even with the SE-ResNet18 backbone which has multiple squeeze-and-excitation attention layers, BoosterSHOT (Tr) improves over MVDeTr in MODA and recall, while the other two metrics are slightly higher. This supports our claim that the attention methods in HAM are able to boost performance by attending to each homography plane, which the backbone is unable to do since it is shared across all homographies.

Table 5. Performance with SE-ResNet18 backbone architecture on Wildtrack

Method	MODA	MODP	precision	recall
MVDeTr	92.2 ± 0.06	81.4 ± 0.21	96.4 ± 0.35	95.7 ± 0.40
BoosterSHOT (Tr)	93.4 ± 0.32	81.5 ± 0.21	96.7 ± 0.10	96.6 ± 0.46

Computational cost, memory consumption and runtime. We analyze the computational cost and memory consumption, along with the runtime of our method and compare with several others to show the benefits of our approach. Through our variable number of selected channels and homographies, we were able to reduce the number of FLOPs and overall runtime without sacrificing performance. More detailed results are outlined in our supplementary.

6. Conclusion

We propose a homography attention module (HAM) as a way to improve across all existing multiview pedestrian detection approaches. HAM consists of a channel gate module that selects the most important channels for each homography and a spatial gate module that applies spatial attention for each homography. We outline an end-to-end multiview pedestrian detection framework (Booster-SHOT) taking insight from previous approaches while also incorporating our proposed module. We report new state-of-the-art performance on standard benchmarks for both Booster-SHOT and previous approaches with HAM and provide extensive empirical evidence that our conjectures and design choices are logically sound. As noted in our supplementary, generalization to novel camera views, bridging the domain gap between synthetic and real-world data, experimentation on large-scale camera models, and extending to tracking and additional analytics for more viability are all areas where further research is needed.

²<https://github.com/huggingface/pytorch-image-models>

References

- [1] Davide Abati, Jakub Tomczak, Tijmen Blankevoort, Simone Calderara, Rita Cucchiara, and Babak Ehteshami Bejnordi. Conditional channel gated networks for task-aware continual learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 6
- [2] Hamid Aghajan and Andrea Cavallaro. *Multi-camera networks: principles and applications*. Academic press, 2009. 1
- [3] Alexandre Alahi, Laurent Jacques, Yannick Boursier, and Pierre Vandergheynst. Sparsity driven people localization with a heterogeneous network of cameras. *Journal of Mathematical Imaging and Vision*, 2011. 1
- [4] Jimmy Ba, Volodymyr Mnih, and Koray Kavukcuoglu. Multiple object recognition with visual attention. *arXiv preprint arXiv:1412.7755*, 2014. 3
- [5] Pierre Baqué, François Fleuret, and Pascal Fua. Deep occlusion reasoning for multi-camera multi-target detection. In *International Conference on Computer Vision (ICCV)*, 2017. 2
- [6] Jerome Berclaz, Francois Fleuret, Engin Turetken, and Pascal Fua. Multiple object tracking using k-shortest paths optimization. *Transactions on pattern analysis and machine intelligence (T-PAMI)*, 2011. 2
- [7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision (ECCV)*, 2020. 1, 3
- [8] Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Timur Bagautdinov, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. Wildtrack: A multi-camera hd dataset for dense unscripted pedestrian detection. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 5
- [9] Tatjana Chavdarova and François Fleuret. Deep multi-camera people detection. In *International Conference on Machine Learning and Applications (ICMLA)*, 2017. 2
- [10] Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In *Conference on computer vision and pattern recognition (CVPR)*, 2017. 3
- [11] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International Conference on Machine Learning (ICML)*, 2020. 3
- [12] Adam Coates and Andrew Y Ng. Multi-camera object detection for robotics. In *International Conference on Robotics and Automation (ICRA)*, 2010. 2
- [13] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *International conference on computer vision (ICCV)*, 2017. 1, 3
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3
- [15] Martin Engilberge, Haixin Shi, Zhiye Wang, and Pascal Fua. Two-level data augmentation for calibrated multi-view detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 128–136, 2023. 2, 3
- [16] Francois Fleuret, Jerome Berclaz, Richard Lengagne, and Pascal Fua. Multi-camera people tracking with a probabilistic occupancy map. *Transactions on pattern analysis and machine intelligence (T-PAMI)*, 2007. 2
- [17] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 3
- [18] Xin Gao, Yijin Xiong, Guoying Zhang, Hui Deng, and Kangkang Kou. Exploiting key points supervision and grouped feature fusion for multiview pedestrian detection. *Pattern Recognition*, 131:108866, 2022. 2, 3
- [19] Zilin Gao, Jiangtao Xie, Qilong Wang, and Peihua Li. Global second-order pooling convolutional networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3
- [20] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 2021. 1
- [21] Meng-Hao Guo, Tian-Xing Xu, Jiang-Jiang Liu, Zheng-Ning Liu, Peng-Tao Jiang, Tai-Jiang Mu, Song-Hai Zhang, Ralph R Martin, Ming-Ming Cheng, and Shi-Min Hu. Attention mechanisms in computer vision: A survey. *arXiv preprint arXiv:2111.07624*, 2021. 1, 3
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on computer vision and pattern recognition (CVPR)*, 2016. 3, 8
- [23] Qibin Hou, Li Zhang, Ming-Ming Cheng, and Jiashi Feng. Strip pooling: Rethinking spatial pooling for scene parsing. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [24] Qibin Hou, Daquan Zhou, and Jiashi Feng. Coordinate attention for efficient mobile network design. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [25] Yunzhong Hou and Liang Zheng. Multiview detection with shadow transformer (and view-coherent data augmentation). In *ACM International Conference on Multimedia*, 2021. 1, 2, 3, 5
- [26] Yunzhong Hou, Liang Zheng, and Stephen Gould. Multi-view detection with feature perspective transformation. In *European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 3, 5, 8
- [27] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Conference on computer vision and pattern recognition (CVPR)*, 2018. 1, 3
- [28] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances in neural information processing systems (NuerIPS)*, 2015. 3

- [29] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *European Conference on Computer Vision (ECCV)*, 2018. 3
- [30] HyunJae Lee, Hyo-Eun Kim, and Hyeonseob Nam. Srm: A style-based recalibration module for convolutional neural networks. In *International Conference on Computer Vision (ICCV)*, 2019. 3
- [31] Wei-Yu Lee, Ljubomir Jovanov, and Wilfried Philips. Multi-view target transformation for pedestrian detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 90–99, 2023. 2
- [32] Jiang-Jiang Liu, Qibin Hou, Ming-Ming Cheng, Changhu Wang, and Jiashi Feng. Improving convolutional networks with self-calibrated convolutions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [33] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. *Advances in neural information processing systems (NeurIPS)*, 2018. 3
- [34] Diganta Misra, TriKay Nalamada, Ajay Uppili Arasanipalai, and Qibin Hou. Rotate to attend: Convolutional triplet attention module. In *Winter Conference on Applications of Computer Vision (WACV)*, 2021. 3
- [35] Volodymyr Mnih, Nicolas Heess, Alex Graves, et al. Recurrent models of visual attention. In *Advances in neural information processing systems (NeurIPS)*, 2014. 3
- [36] Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa, Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention u-net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*, 2018. 3
- [37] Jongchan Park, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Bam: Bottleneck attention module. *The British Machine Vision Conference (BMVC)*, 2018. 3
- [38] Zequn Qin, Pengyi Zhang, Fei Wu, and Xi Li. Fcanet: Frequency channel attention networks. In *International Conference on Computer Vision (ICCV)*, 2021. 3
- [39] Gemma Roig, Xavier Boix, Horesh Ben Shitrit, and Pascal Fua. Conditional random fields for multi-camera object detection. In *International Conference on Computer Vision (ICCV)*, 2011. 2
- [40] Abhijit Guha Roy, Nassir Navab, and Christian Wachinger. Recalibrating fully convolutional networks with spatial and channel “squeeze and excitation” blocks. *Transactions on medical imaging*, 2018. 3
- [41] Aswin C. Sankaranarayanan, Ashok Veeraraghavan, and Rama Chellappa. Object detection, tracking and recognition for multiple smart cameras. *Proceedings of the IEEE*, 2008. 1, 2
- [42] Liangchen Song, Jialian Wu, Ming Yang, Qian Zhang, Yuan Li, and Junsong Yuan. Stacked homography transformations for multi-view pedestrian detection. In *International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 3, 4, 5, 6, 8
- [43] Xiaoxiao Sun and Liang Zheng. Dissecting person re-identification from the viewpoint of viewpoint. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 5
- [44] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *Conference on computer vision and pattern recognition (CVPR)*, 2017. 3
- [45] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Computer vision and pattern recognition (CVPR)*, 2018. 3
- [46] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *European conference on computer vision (ECCV)*, 2018. 1, 3, 6
- [47] Saining Xie, Sainan Liu, Zeyu Chen, and Zhuowen Tu. Attentional shapecontextnet for point cloud recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [48] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning (ICML)*, 2015. 3
- [49] Zongxin Yang, Linchao Zhu, Yu Wu, and Yi Yang. Gated channel transformation for visual recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [50] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 8
- [51] Yuhui Yuan, Lang Huang, Jianyuan Guo, Chao Zhang, Xilin Chen, and Jingdong Wang. Ocnet: Object context for semantic segmentation. *International Journal of Computer Vision*, 2021. 1
- [52] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3
- [53] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3