# Wide-Area Crowd Counting: Multi-View Fusion Networks for Counting in Large Scenes
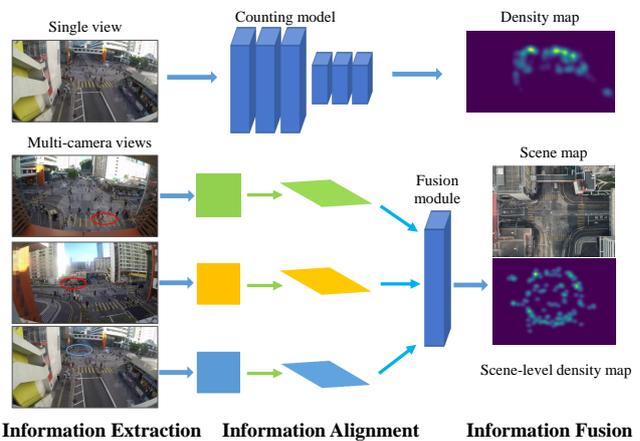
**Qi Zhang · Antoni B. Chan**

**Abstract** Crowd counting in single-view images has achieved outstanding performance on existing counting datasets. However, single-view counting is not applicable to large and wide scenes (*e.g.*, public parks, long subway platforms, or event spaces) because a single camera cannot capture the whole scene in adequate detail for counting, *e.g.*, when the scene is too large to fit into the field-of-view of the camera, too long so that the resolution is too low on faraway crowds, or when there are too many large objects that occlude large portions of the crowd. Therefore, to solve the wide-area counting task requires multiple cameras with overlapping fields-of-view. In this paper, we propose a deep neural network framework for multi-view crowd counting, which fuses information from multiple camera views to predict a scene-level density map on the ground-plane of the 3D world. We consider three versions of the fusion framework: the late fusion model fuses camera-view density map; the naïve early fusion model fuses camera-view feature maps; and the multi-view multi-scale early fusion model ensures that features aligned to the same ground-plane point have consistent scales. A rotation selection module further ensures consistent rotation alignment of the features. We test our 3 fusion models on 3 multi-view counting datasets, PETS2009, DukeMTMC, and a newly collected multi-view counting dataset containing a crowded street intersection. Our methods achieve state-of-the-art results compared to other multi-view counting baselines.

Qi Zhang
City University of Hong Kong
E-mail: qzhang364-c@my.cityu.edu.hk

Antoni B. Chan
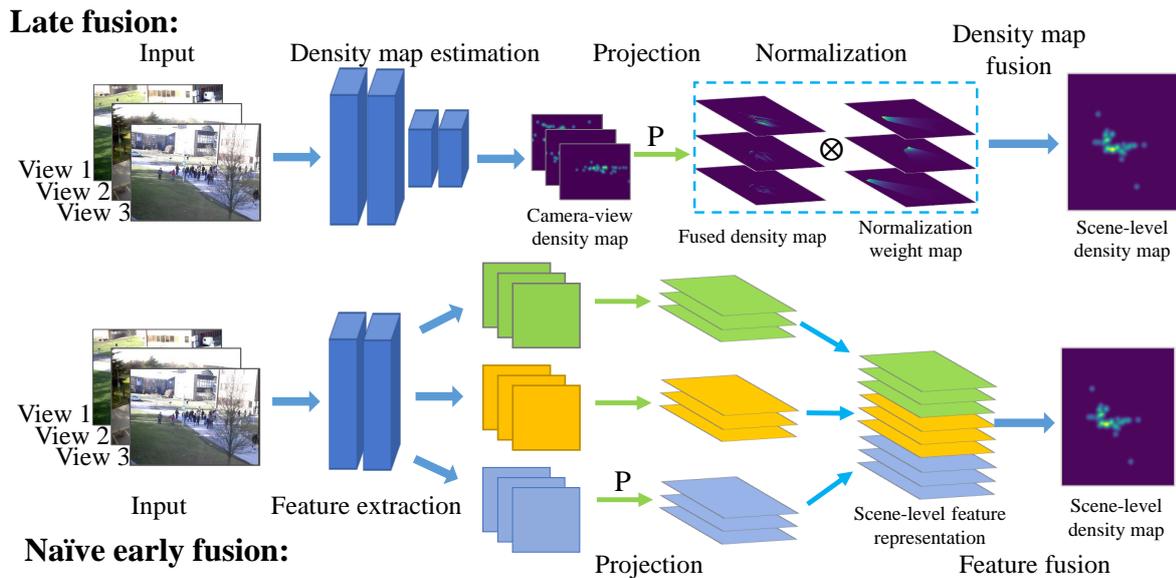City University of Hong Kong
E-mail: abchan@cityu.edu.hk



**Fig. 1** The pipeline of the proposed multi-view fusion framework comparing with the single image counting framework. In the multi-view fusion model, feature maps are extracted from multiple camera views, aligned on the ground-plane, and fused to obtain the scene-level ground-plane density map. The scene map is shown for reference. For single image counting, many people are occluded (in red circles) and in low resolution (in blue circles), which decreases the counting performance.

## 1 Introduction

Crowd counting aims to estimate the number of the people in images or videos. It has a wide range of real-world applications, such as crowd management, public safety, traffic monitoring or urban planning (Sindagi and Patel 2018). For example, crowd counting can detect overcrowding on the railway platform and help with the train schedule planning. Furthermore, the estimated crowd density map provides spatial information of the crowd, which can benefit other tasks, such as hu-

**Fig. 2** The pipeline of our late fusion model and naïve early fusion model for multi-view counting. In the late fusion model, single-view density maps are fused. In the naïve early fusion model, single-view feature maps are fused.

man detection (Eiselein et al. 2013; Kang et al. 2018; Ma et al. 2015) and tracking (Kang et al. 2018; Ren et al. 2018; Rodriguez et al. 2011).

Recently, with the strong learning ability of deep neural networks (DNNs), density map based crowd counting methods have achieved outstanding performance on the existing counting datasets (Cao et al. 2018; Idrees et al. 2018; Sindagi and Patel 2017), where the goal is to count the crowd in a single image. However, a single image view is not adequate to cover a *large* and *wide* scene, such as a large park or a long train platform. For these wide-area scenes, a single camera view cannot capture the whole scene in adequate detail for counting, either because the scene is too large (wide) to fit within the field-of-view of the camera, or the scene is too long so that the resolution is too low in faraway regions. Furthermore, a single view cannot count regions that are still within the scene, but are totally occluded by large objects (*e.g.*, trees, large vehicles, building structures). Therefore, to solve the wide-area counting task requires multiple camera views with overlapping field-of-views, which combined can cover the whole scene and can see around occlusions. The goal of wide-area counting is then to use multiple camera views to estimate the crowd count of the whole scene.

Existing multi-view counting methods rely on foreground extraction techniques and hand-crafted features. Their crowd counting performance is limited by the effectiveness of the foreground extraction, as well as the representation ability of hand-crafted features. Considering the strong learning power of DNNs as well as the

performance progress of single view counting methods using density maps, the feasibility of end-to-end DNNs-based multi-view counting methods should be explored.

In this paper, we propose a DNNs-based multi-view counting method that extracts information from each camera view and then fuses them together to estimate a scene-level ground-plane density map (see Fig. 1). The method consists of 3 stages: 1) *Information extraction* – single view feature maps are extracted from each camera image with DNNs; 2) *Information alignment* – using the camera geometry, the feature maps from all cameras are projected onto the ground-plane in the 3D world so that the same person's features are approximately aligned across multiple views, and properly normalized to remove projection effects; 3) *Information fusion* – the aligned single-view projected feature maps are fused together and used to predict the scene-level ground-plane density map.

As single-view crowd counting is relatively mature, it has well-studied feature extractor and decoder architectures for predicting camera-level density maps. Building from this, multi-view crowd counting can leverage single-view feature extractors and predicted density maps for each camera-view. The key issue then is *what* and *how* to fuse the information from the various cameras into a ground-plane representation for decoding into a ground-plane density map. We consider three variants of fusion to provide a thorough study on the fusion architecture for multi-view counting DNNs models. These three variants differ in *what information is fused* (i.e., single-view density maps or feature

maps) and *how fusion occurs* (i.e., simple concatenation or scale/rotation-aware concatenation). Specifically, the three variants are: 1) concatenation of single-view density maps (denoted as late fusion); 2) concatenation of single-view feature maps (naïve early fusion); 3) scale-aware and rotation-aware concatenation of feature maps (multi-view multi-scale, MVMS/MVMSR).

Specifically, first, in our late-fusion model (see Fig. 2 top), view-level density maps are predicted for each camera view, projected to the ground-plane, and fused for estimating the scene-level density map. This model fuses count-level information, similar to traditional count-based methods (Dittrich et al. 2017; Li et al. 2012; Ma et al. 2012; Maddalena et al. 2014). We also propose a post-projection normalization method that removes the projection effect that distorts the sum of the density maps (and thus the count). Second, in our naïve early fusion model (see Fig. 2 bottom), convolutional feature maps are extracted from each camera view, projected to the ground-plane and fused to predict the scene-level density map. Third, to handle the scale variations of the same person across camera views, our multi-view multi-scale (MVMS) early fusion model (see Fig. 6) extracts features with consistent scale across corresponding locations in the camera views before applying projection and fusion. We consider 2 approaches for selecting the suitable scales, based on distances computed from the camera geometry. To further improve the multi-view fusion performance, a rotation selection module is added in the multi-view fusion step (denoted as MVMSR).

The existing multi-view datasets that can be used for multi-view counting are PETS2009 (Ferryman and Shahrokni 2009) and DukeMTMC (Ristani et al. 2016). However, PETS2009 is not a wide-area scene as it focuses on one walkway, while DukeMTMC is a wide-area scene but does not contain large crowds. To address these shortcomings, we collect a new wide-area dataset from a busy street intersection, which contains large crowds, more occlusion patterns (*e.g.*, buses and cars), and large scale variations. This new dataset more effectively tests multi-view crowd counting in a real-world scene.

In summary, our main contributions are:

1. We propose an end-to-end trainable DNNs-based multi-view crowd counting framework, which fuses information from multiple camera views to obtain a scene-level density map.
2. We propose 3 fusion models based on our multi-view framework (late fusion, naïve early fusion, and multi-view multi-scale early fusion), which achieve better counting accuracy compared to baselines.
3. We propose a rotation selection module based on rotation equivariant networks to further improve the multi-view fusion by considering the geometric properties of the average-height projection.
4. We collect a real-world wide-area counting dataset consisting of multiple camera views, which will advance research on multi-view wide-area counting.

The remainder of this paper is organized as follows. In Section 2, existing single-view and multi-view counting methods are reviewed, and the rotation neural networks are introduced. In Section 3, the proposed two DNNs-based multi-view counting models (both late fusion and naïve early fusion model) are presented. In Section 4, the multi-view multi-scale early fusion model with scale selection and rotation selection module is presented. In Section 5, we conduct experiments on multi-view counting datasets.

## 2 Related Work

In this section, we review methods for crowd counting from single-view and multi-view cameras, as well as rotation equivariant/invariant networks.

### 2.1 Single-view counting

*Traditional methods.* Traditional single-view counting methods can be divided into 3 categories (Chen et al. 2013; Sindagi and Patel 2018): detection, regression, and density map methods. Detection methods try to detect each person in the images by extracting hand-crafted features (Viola and Jones 2004; Sabzmeydani and Mori 2007; Wu and Nevatia 2007) and then training a classifier (Joachims 1998; Viola et al. 2005; Gall et al. 2011) using the extracted features. However, the detection methods do not perform well when the people are heavily occluded, which limits their application scenarios. Regression methods extract image features (Chan et al. 2008; Cheng et al. 2014; Junior et al. 2010; Krizhevsky et al. 2012) and learn a mapping directly to the crowd count (Chan and Vasconcelos 2012; Chen et al. 2012; Paragios and Ramesh 2001; Marana et al. 1998). However, their performance is limited by the weak representation power of the hand-crafted low-level features. Instead of directly obtaining the counting number, Lempitsky and Zisserman (2010) proposed to estimate density maps, where each pixel in the image contains the local crowd density, and the count is obtained by summing over the density map. Traditional density map methods learn the mapping between the hand-crafted local features and the density maps (Lempitsky and Zisserman 2010; Pham et al. 2015; Wang and Zou 2016; Xu and Qiu 2016).

*DNNs-based methods.* DNNs-based crowd counting has mainly focused on density map estimation. The first networks used a standard CNN (Zhang et al. 2015) to directly estimate the density map from an image. Scale variation is a critical issue in crowd counting, due to perspective effects in the image (Yan et al. 2019a; Liu et al. 2019b; Xu et al. 2019). Zhang et al. (2016) proposed the multi-column CNN (MCNN) consisting of 3 columns of different receptive field sizes, which can model people of different scales. Sam et al. (2017) added a switching module in the MCNN structure to choose the optimal column to match the scale of each patch. Onoro-Rubio and López-Sastre (2016) proposed to use the patch pyramid as input to extract multi-scale features. Similarly, Kang and Chan (2018) used an image pyramid with a scale-selecting attention block to adaptively fuse predictions on different scales.

Recently, more sophisticated network structures have been proposed and extra information is explored to advance the counting performance (Shi et al. 2018; Idrees et al. 2018; Wang et al. 2019; Ranjan et al. 2018; Cao et al. 2018; Li et al. 2018; Liu et al. 2018; Shen et al. 2018; Jiang et al. 2019; Liu et al. 2019c). Sindagi and Patel (2017) incorporated global and local context information in the crowd counting framework, and proposed the contextual pyramid CNN (CP-CNN). Idrees et al. (2018) proposed composition loss, implemented through multiple dense blocks after branching off the base networks. Li et al. (2018) replaced pooling operations in the CNN layers with dilated kernels to deliver larger reception fields and achieved better counting performance. Kang et al. (2017) proposed an adaptive convolution neural network (ACNN) that uses side information (camera angle and height) to include context into the counting framework. Many methods have focused on the perspective change issue in the counting task. Cao et al. (2018) extracted multi-scale features with a scale aggregation module and generated high-resolution density maps by using a set of transposed convolutions. Shi et al. (2019) proposed to estimate the perspective map and use it to adaptively fuse the multi-scale output density maps. Yan et al. (2019b) proposed perspective-guided convolution (PGC) to utilize perspective information instead of multi-scale or multi-column architectures. Yang et al. (2020) proposed to estimate a perspective factor to warp the input images to correct the perspective distortions. Lian et al. (2019) proposed a regression guided detection network (RD-Net) for RGB-D crowd counting. Liu et al. (2019a) proposed Recurrent Attentive Zooming Network to zoom high density regions for higher-precision counting and localization.

All these methods are using DNNs to estimate a density map on the image plane of a single camera-view, with different architectures improving the performance across scenes and views. In contrast, in this paper, we focus on fusing multiple camera views of the same scene to obtain a ground-plane density map in the 3D world. These single-view methods serve as the backbone single-view feature extractors for our multi-view fusion networks.

## 2.2 Multi-view counting

Existing multi-view counting methods can be divided into 3 categories: detection/tracking, regression, and 3D cylinder methods. The detection/tracking methods first perform detection or tracking on each scene and obtain single-view detection results. Then, the detection results from each view are integrated by projecting the single-view results to a common coordinate system, *e.g.*, the ground plane or a reference view. The count of the scene is obtained by solving a correspondence problem (Dittrich et al. 2017; Li et al. 2012; Ma et al. 2012; Maddalena et al. 2014). Regression based methods first extract foreground segments from each view, then build the mapping relationship of the segments and the count number with a regression model (Ryan et al. 2014; Tang et al. 2014). 3D cylinder-based methods try to find the people's locations in the 3D scene by minimizing the gap between the people's 3D positions projected into the camera view and the single view detection (Ge and Collins 2010).

These multi-view counting methods are mainly based on hand-crafted low-level features and regression or detection/tracking frameworks. Regression-based methods only give the global count, while detection/tracking methods cannot cope well with occlusions when the scene is very crowded. In contrast to these works, our approach is based on predicting the ground-plane density map in the 3D world by fusing the information across camera views using DNNs. Two advantages of our approach are the abilities to learn the feature extractors and fusion stage in end-to-end training, and to estimate the spatial arrangement of the crowd on the ground plane. While the previous methods are mainly tested on PETS2009, which only contains low/moderate crowd numbers on a walkway, here we test on a newly collected dataset comprising a real-world scene of a street intersection with large crowd numbers, vehicles, and occlusions.

A preliminary conference version of this work appears in Zhang and Chan (2019). This journal version contains the following extensions: 1) more details about

the scale selection module are added, specifically the rationale of the multi-view scale selection guided by the distance map; 2) a new rotation selection module is proposed to consider the stretching effect of the fixed average-height projection, which further boosts the performance compared to the model in Zhang and Chan (2019); 3) more experiments and ablation studies, including more evaluation metrics (MAE, MSE, NAE and GAME), experiments of new and updated comparison methods ('feature concatenation', 'stitching', and updated 'Detection+ReID' method), experiments showing how multi-cameras improve single-view counting performance for each multi-camera counting methods, more ablation studies on the rotation modules (filter number, layer number and quantization angle), experiments with more recent backbone networks, comparison results with different module settings and methods and on another test set in DukeMTMC, and the running speed comparison of different methods.

Finally, following our conference paper (Zhang and Chan 2019), a subsequent work (Zheng et al. 2021) enhances the late fusion model's performance by modeling the correlation between each pair of views for cross view fusion.

## 2.3 Rotation equivariant/invariant networks

Rotation equivariance or invariance relates to the DNNs' robustness to rotation changes of the input image. Rotation equivariance means the output is accordingly rotated if the input is rotated, which is useful for the dense prediction tasks, like semantic segmentation or density map estimations. Rotation invariance means the output is invariant no matter how the input is rotated, which is useful for classification tasks.

To enhance the networks' robustness to rotations, the easiest method is to use the data augmentation, namely rotating the original examples multiple times and training the network on the rotated versions. Jaderberg et al. (2015) introduced the Spatial Transformer which can spatially manipulate data within the network, giving neural networks the ability to actively spatially transform feature maps. Rotation equivariant and invariant networks have also been proposed to improve the rotation robustness. Laptev et al. (2016) uses multiple rotated examples as inputs into a shared network to extract features at multiple rotations, and then uses a max-pooling layer among these features to obtain rotation-invariant features. In addition to image rotating, feature maps can also be rotated. For example, Dieleman et al. (2016) and Cohen and Welling (2016) obtained rotation robustness by rotating the feature maps 3 times, by 90 degree each time, and then used

average or max-pooling operations. Besides images and feature maps, rotation robustness can also be obtained by rotating the kernel/filter. Gao and Ji (2017) provided an example of how to use rotated kernels, but the weakness of their method is that the kernel size is fixed (3*3), and rotation angle is limited (45 degree each time, not arbitrary). Marcos et al. (2017) performed arbitrary rotations and the orientation pooling was used instead of max pooling to get the rotation-equivariance. Weiler et al. (2018) proposed the steerable filter CNNs, which employed steerable filters to compute orientation dependent responses without suffering interpolation artifacts from filter rotation, and used group convolutions for an equivariant mapping. Recently, rotation equivariant/invariant networks have been utilized in 3D recognition tasks, such as CubeNet (Worrall and Brostow 2018), ClusterNet (Chen et al. 2019).

In contrast to these methods that aim to obtain robustness to rotations, we use the rotation equivariant/invariant networks to negate the effects of the projection operation from the camera-view to the ground-plane. Specifically, the multi-rotated filters are used to reduce the influence of the average-height projection on the extracted features, which improves the multi-view fusion counting performance.

## 3 Multi-View Counting via Multi-View Fusion

For multi-view counting, we assume that the cameras are fixed, the camera calibration parameters (both intrinsic and extrinsic) are known, and that the camera frames across views are synchronized. Given the set of multi-view images, the goal is to predict a scene-level density map defined on the ground-plane of the 3D scene (see Fig. 1). The ground-truth ground-plane density map is obtained in a similar way as the traditional camera-view density map – the ground-plane annotation map is obtained using the ground-truth 3D coordinates of the people, which is then convolved by a fixed-width Gaussian kernel to obtain the density map on the ground-plane.

In the following two sections, we propose three fusion approaches for multi-view counting: 1) the *late fusion* model projects camera-view density maps onto the ground plane and then fuses them together, and requires a projection normalization step; 2) the *naïve early fusion* model projects camera-view feature maps onto the ground plane then fuses them; 3) to handle inter-view and intra-view scale variations, the *multi-view multi-scale early fusion* model (MVMS) selects features scales to be consistent across views when projecting to the same ground-plane point, and uses rotation

| FCN-7 | |
|---|---|
| Layer | Filter |
| conv 1 | $16 \times 1 \times 5 \times 5$ |
| conv 2 | $16 \times 16 \times 5 \times 5$ |
| pooling | $2 \times 2$ |
| conv 3 | $32 \times 16 \times 5 \times 5$ |
| conv 4 | $32 \times 32 \times 5 \times 5$ |
| pooling | $2 \times 2$ |
| conv 5 | $64 \times 32 \times 5 \times 5$ |
| conv 6 | $32 \times 64 \times 5 \times 5$ |
| conv 7 | $1 \times 32 \times 5 \times 5$ |

| Fusion | |
|---|---|
| Layer | Filter |
| concat | - |
| conv 1 | $64 \times n \times 5 \times 5$ |
| conv 2 | $32 \times 64 \times 5 \times 5$ |
| conv 3 | $1 \times 32 \times 5 \times 5$ |

**Table 1** FCN-7 backbone and fusion module. The Filter dimensions are output channels, input channels, and filter size ($w \times h$).
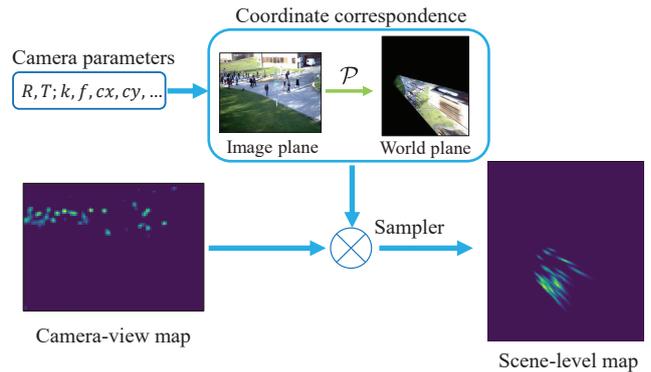
selection to handle rotation effects of the projection (MVMSR).

We note that there are differences between fusing density maps and fusing feature maps, since they contain different information. Density maps only contain location information of the crowd, but do not contain identity information of people (person appearance features). Thus, fusing density maps requires combining / aligning local density information across projected views, but may suffer from errors due to ambiguity or distortion caused by the 2D to 3D projection. In contrast, feature maps contain identity (appearance) information that can help to find correspondences of the same person in different views on the ground-plane. However, since the person is a different distance from each camera, this information is present at different scales among the camera views, which makes learning the correspondences more difficult because the DNN should see all combinations of scales to become scale invariant. Thus to alleviate the issue of scale variations among cameras, we propose a scale-aware fusion step, which uses image pyramids to select features at the same scale before projection. In this way, all the features projected onto the ground-plane are at the same scale, and the relationships among features is easier to learn.

We first present the common components, and then the 3 fusion models.

### 3.1 Backbone for camera views

A fully-convolutional network (denoted as FCN-7) is used on each camera view to extract image feature maps or estimate a corresponding view-level density map. The FCN-7 settings are shown in Table 1. For the ablation study, CSR-Net (Li et al. 2018) and LCC (Liu et al. 2020) are also used as feature backbone (see Section 6.3.1).



**Fig. 3** The projection module to transform camera-view maps to a ground-plane representation. Here the camera-view map is visualized as a density map.
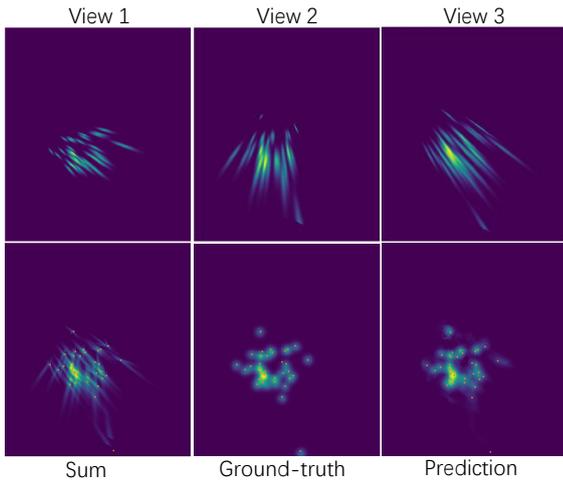
### 3.2 Image to ground-plane projection

As we assume that the intrinsic and extrinsic parameters of the cameras are known, the projection from a camera's 2D image space to a 3D ground-plane representation can be implemented as a differentiable fixed-transformation module (see Fig. 3). The 3D height (z-coordinate) corresponding to each image pixel is unknown. Since the view-level density maps are based on head annotations and the head is typically visible even during partial occlusion, we assume that each pixel's height in the 3D world is a person's average height (1750 mm). The camera parameters together with the height assumption are used to calculate the correspondence mapping $\mathcal{P}$ between 2D image coordinates and the 3D coordinates on the 3D average-height plane. Finally, the Sampler from the Spatial Transformer Networks (Jaderberg et al. 2015) is used to implement the projection, resulting in the ground-plane representation of the input map.

### 3.3 Late fusion model

and then project them to the ground-plane for fusion and obtaining the scene-level density map, where the intersections of the projected Gaussians are close to the people locations on the ground-plane (see Fig. 4) and are mapped to the Gaussian kernels of the ground-plane density map.

The main idea of the late fusion model is to first estimate the crowd density maps in each camera view, and then project them to the ground-plane for fusion and prediction of the scene-level density map. As shown in Fig. 4, the intersections of the projected Gaussians will be close to the people locations on the ground-plane, and thus the fusion network aims to transform the inter-
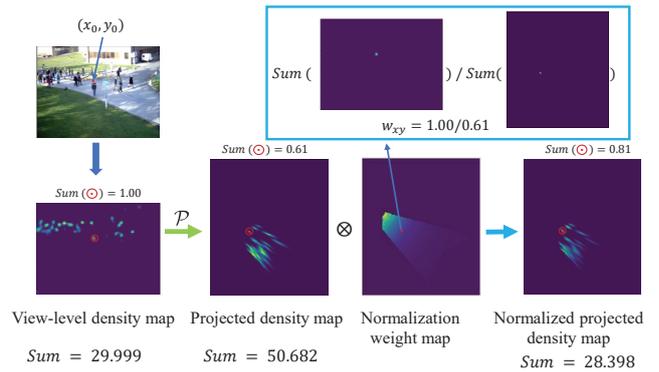
**Fig. 4** Example of single-view density maps projected on to the ground-plane and their summation. The intersections of the projected Gaussians are close to the people locations on the ground-plane. The orange dots are the ground-truth annotations.



**Fig. 5** The projection normalization process for the late fusion model. $Sum$ is the sum of the whole density map, while $Sum(\odot)$ is the sum over the circled region (diameter 5).

section points into the Gaussian kernels on the ground-plane. In particular, the late fusion model consists of 3 stages (see Fig. 2 top): 1) estimating the camera-view density maps using FCN-7 on each view; 2) projecting the density maps to the ground-plane representation using the projection module; 3) concatenating the projected density maps channel-wise and then applying the Fusion module to obtain the scene-level density map. The network settings for the fusion network are presented in Table 1.

*Projection Normalization.* One problem is that the density map is stretched during the projection step, and thus the sum of the density map changes after the projection. Considering that the density map is composed of a sum of Gaussian kernels, each Gaussian is stretched differently depending on its location in the image plane. To address this problem, we propose a normalization method to ensure that the sum of each Gaussian kernel remains the same after projection (see Fig. 5). In particular, let $(x_0, y_0)$ and $(x, y)$ be the corresponding points in the image plane and the 3D world ground-plane representation. The normalization weight $w_{xy}$ for ground-plane position $(x, y)$ is

$$w_{xy} = \frac{\sum_{i,j} D_{x_0,y_0}(i,j)}{\sum_{m,n} \mathcal{P}(D_{x_0,y_0}(m,n))}, \tag{1}$$

where $D_{x_0,y_0}$ denotes an image-space density map containing **only one Gaussian kernel centered** at $(x_0, y_0)$, $\mathcal{P}$ is the projection operation from image space to ground plane representation, the summation operation is over the whole camera view map or projected ground-plane map, and $(i, j)$ and $(m, n)$ are the image coordinates
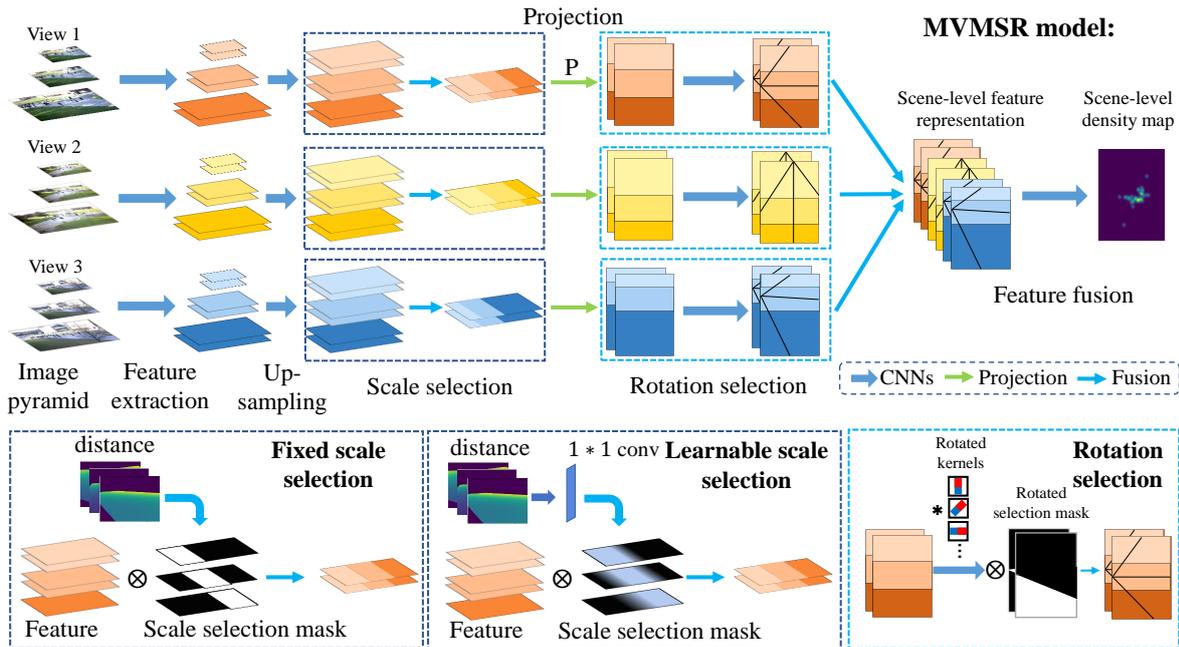
and ground-plane coordinates, respectively. The normalization map $W = [w_{xy}]$ for each camera is element-wise multiplied to the corresponding projected density map before concatenation. As illustrated in Fig. 5, to visualize the effect of the projection normalization, we choose a small circular region (diameter 5), and calculate its sum before and after using the normalization. After applying the projection normalization, the local ground-plane sum is more consistent with the corresponding image-based sum. Likewise, the sums of the whole ground-plane and the whole image are also more consistent after applying the normalization.

### 3.4 Naïve early fusion model

The naïve early fusion model directly fuses the feature maps from all the camera-views to estimate the ground-plane density map. Similar to the late fusion model, we implement the early fusion model by replacing the density map-level fusion with feature-level fusion (see Fig. 2 bottom). Specifically, the naïve early fusion model consists of 3 stages: 1) extracting feature maps from each camera view using the first 4 convolution layers of FCN-7; 2) projecting the image feature maps to the ground-plane representation using the projection module; 3) concatenating the projected feature maps and applying the Fusion module to estimate the scene-level density map. Note that the projection normalization step used in the late fusion model is not required for the early fusion model, since feature maps do not have the same interpretation of summation yielding a count.

## 4 Multi-view multi-scale early fusion model

Intra-view scale variations are an important issue in single-view counting, as people will appear with differ-

**Fig. 6** The pipeline of multi-view multi-scale early fusion model (MVMS) with rotation selection module (MVMSR). First, multi-scale feature maps are extracted with an image pyramid. The multi-scale feature maps are up-sampled to the same size. The scale selection module (the dotted box) ensures the scales of features that represent the same ground-plane point are consistent across all views. The scale-consistent features are projected to the average-height plane and then fused to obtain the scene-level density map. Two kinds of scale selection strategies (the two dotted boxes on the right) are utilized: the fixed scale selection uses the distance information relative to a reference distance, and learnable scale selection makes the reference distance a learnable parameter. For MVMSR, a rotation selection module is added after the projection step and before the fusion step in order to remove mis-aligned rotations caused by the projection step.

ent sizes in the image due to perspective effects. Using multiple views increases the severity of the scale variation issue; in addition to intra-view scale variation, multi-view images have inter-view scale variations, where the same person will appear at different scales across multiple views. This inter-view scale variation may cause problems during the fusion stage as there are a combinatorial number of possible scales appearing across all views, which the network needs to be invariant to. To address this problem, we extract feature maps at multiple scales, and then perform scale selection so that the projected features are at consistent scales across all views (*i.e.*, a given person's features are at the same scale across all views).
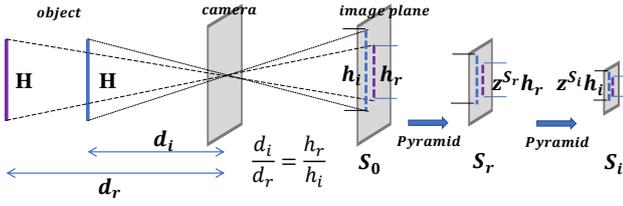
Our proposed multi-view multi-scale (MVMS) early fusion architecture is shown in Fig. 6. The MVMS fusion model consists of 4 stages: 1) extracting multi-scale feature maps by applying the first 4 convolution layers of FCN-7 on an image pyramid for each camera view; 2) upsampling all the feature maps to the largest size, and then selecting the scales for each pixel in each camera-view according to the scene geometry; 3) projecting the scale-consistent feature maps to the ground-plane representation using the projection module; 4) fusing the

projected features and predicting a scene-level density map using the fusion module. We consider 2 strategies for selecting the consistent scales, fixed scale selection and learnable scale selection.

To further boost the multi-view fusion process, a rotation selection module is added before the fusion step in the MVMS model, denoted as *MVMSR*. In the rotation selection module, the projected feature maps are convolved with multiple rotated versions of the filters, and then combined by selecting among the rotated features based on the camera geometry.

### 4.1 Scale selection module

In the camera pinhole model, an object's scale in an image is influenced by the object's distance to the camera (see Fig. 7). Therefore, the distance-to-camera information can be used to select the scale in an image pyramid to achieve scale consistency across multiple views. A distance map of each view can be calculated from the camera extrinsic parameters and the average person height. The projection operation $\mathcal{P}(x_0, y_0, h_{avg})$ is the projection of the image view coordinate $(x_0, y_0)$ to the 3D world coordinates on the average height plane

**Fig. 7** The relationship between distance-to-camera and object scale in a camera pinhole model.



**Fig. 8** Visualization of camera distance maps of PETS2009.



**Fig. 9** The fixed and learnable scale selection masks for PETS2009.

$h_{avg}$. Then the distance-to-camera $d(x_0, y_0)$ (see Fig. 8) is calculated by transforming to the camera coordinate system, where the camera center is the origin,

$$d(x_0, y_0) = ||R\mathcal{P}(x_0, y_0, h_{avg}) + T||, \tag{2}$$

where $R$ and $T$ are the camera extrinsic parameters, rotation matrix and translation, respectively.

Next we show how to use the distance information to compute the scale according to the pinhole camera model. Consider an image pyramid with zoom factor $z$ between neighboring scales. Let $H$ be the height of the object in the 3D world ($H = h_{avg}$ here), and define $h_r$ as the height of the object on the image (at scale 0) when the object is at a reference distance $d_r$ from the camera. The same object appears in the image with a different height $h_i$ when it is at distance $d_i$ from the camera. According to camera pinhole model, we have $d_i/d_r = h_r/h_i$. In the image pyramid, the object's height is $z^{S_r}h_r$ in image scale $S_r$ (at distance $d_r$) and $z^{S_i}h_i$ in image scale $S_i$ (at distance $d_i$). Thus, to achieve scale consistency, where the heights are equal in the selected image scales, we require that
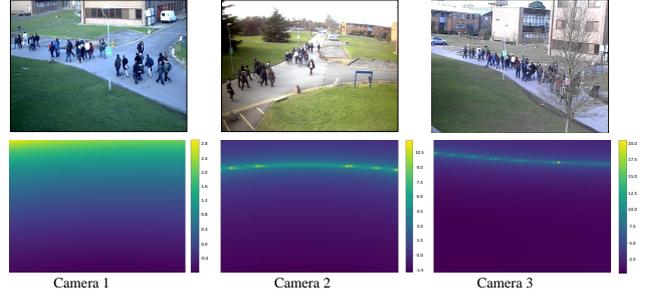
$$z^{S_r}h_r = z^{S_i}h_i. \tag{3}$$

Solving for $S_i$, we obtain the scale required for the object at distance $d_i$ to be consistent with the object at reference distance $d_r$ and at reference scale $S_r$,
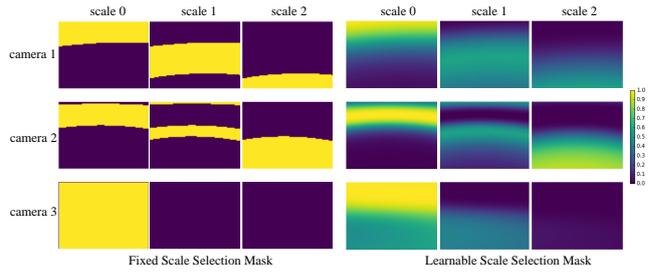
$$S_i = S_r - \log_z(d_i/d_r). \tag{4}$$

### 4.1.1 Fixed scale-selection

The fixed scale selection strategy is illustrated in Fig. 6 (bottom-left). For a given camera, let $\{F_0, \cdots, F_n\}$ be the set of feature maps extracted from the image pyramid, and then upsampled to the same size. Here $F_0$ is the original scale and $F_n$ is the smallest scale. A distance map is computed according to (2) for the camera-view, where $d(x_0, y_0)$ is the distance between the camera's 3D location and the projection of the point $(x_0, y_0)$ into the 3D-world (on the average height plane). A scale selection map $S$, where each value corresponds to the

selected scale for that pixel, is computed according to (4),

$$S(x_0, y_0) = S_r - \lfloor \log_z \frac{d(x_0, y_0)}{d_r} \rfloor, \tag{5}$$

where $\lfloor \cdot \rfloor$ is the floor function. $d_r$ and $S_r$ are the reference distance and the corresponding reference scale number, which are the same for all camera-views. In our experiments, we set the reference distance $d_r$ as the distance value for the center pixel of the first view, and $S_r$ as the middle scale of the image pyramid. Given the scale selection map $S$, the feature maps across scales are merged into a single feature map, $F = \sum_i \mathbb{1}(S = i) \otimes F_i$, where $\otimes$ is element-wise multiplication, and $\mathbb{1}$ is an element-wise indicator function.

### 4.1.2 Learnable scale-selection

The fixed scale selection strategy requires setting the reference distance and reference scale parameters. To make the scale selection process more adaptive to the view context, a learnable scale-selection model is considered (see Fig. 6 (bottom-right)),

$$S(x_0, y_0) = b + k \log_z \frac{d(x_0, y_0)}{d_r}, \tag{6}$$

where the learnable parameter $b$ corresponds to the reference scale, and $k$ adjusts the reference distance. The learnable scale selection can be implemented as a $1 \times 1$

**Fig. 10** Rotation distortion of projected features caused by the average-height projection (using the original image for visualization).



**Fig. 11** The rotation selection layer. The same kernel is padded, rotated and then convolved with the projected features. The rotation selection mask is used to select and fuse the multi-rotated features.

convolution on the log distance map. Then, a soft scale selection mask $M_i$ for scale $i$ can be obtained,

$$M_i(x_0, y_0) = \frac{e^{-(S(x_0,y_0)-i)^2}}{\sum_{j=0}^{n} e^{-(S(x_0,y_0)-j)^2}}. \tag{7}$$

Note that $M_i(x_0, y_0)$ has values between 0 and 1, and sums to 1 across scales, $\sum_{i=0}^{k} M_i(x_0, y_0) = 1$. The scale consistent feature map is then

$$F = \sum_{i=0}^{k} M_i \otimes F_i, \tag{8}$$

which is equivalent to a per-pixel soft-attention mechanism across scales, where $M_i$ is the pixel-wise soft-attention mask.

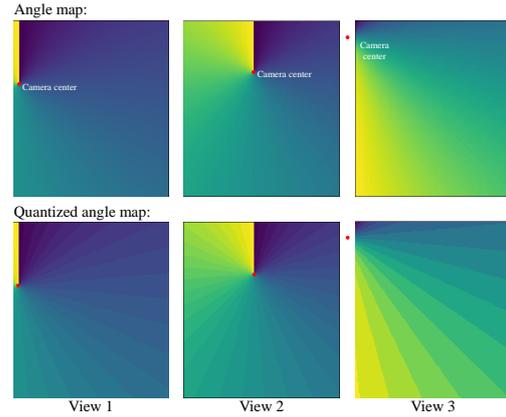The fixed and learnable scale selection masks for PETS2009 are shown in Fig. 9. The fixed scale selection produces binary masks and learnable scale selection produces soft masks. The learnable scale selection gives more freedom for the networks to fuse the multi-scale features, especially on the edges of each scale's masks.

### 4.2 Rotation selection module

In the projection module, the average-height assumption is utilized in which all pixels in each view are assumed to have the average height. The average-height



**Fig. 12** Examples of the rotation selection masks for PETS2009. The first row shows the angle maps for each view, which show the rotation angle for each pixel. Brighter color means larger rotation angles. The second row shows the quantized angle maps, where a finite number of rotation angles are used to reduce computation complexity.

projection is applicable due to the head annotations in the datasets. On the other hand, the average height projection makes the feature patterns stretched along the view ray. Therefore, the features are "rotated" to the view ray direction after the projection (see Fig. 10). To further improve the multi-view fusion, a rotation selection module is proposed and used before the multi-view fusion step in order to counteract this phenomenon.

The rotation selection layer is illustrated in Fig. 11, where "tall" aspect-ratio kernels ($k_1$ by $k_2$) are adopted due to the stretched patterns. First, the kernel is padded to be square to ensure the feasibility of the arbitrary rotation of the kernel, where the square size is $\lceil \sqrt{2} \max(k_1, k_2) \rceil$. Second, the kernels are rotated over a range of angles $\{r_0, r_1, \cdots, r_m\}$ (decided by the rotation selection map, see next paragraph), which can be implemented using the Sampler from (Jaderberg et al. 2015). Third, each rotated kernel is convolved with the projected feature maps, resulting in multi-rotated features $\{F_0, ..., F_m\}$. Finally, the multi-rotated features are selected and fused with rotation selection masks.

The rotation selection mask is calculated with the camera parameters on the average-height plane. Suppose $W_c = (x_c, y_c)$ are the coordinates on the scene-level plane (after projection), and the $(R, T)$ are the camera extrinsic parameters. Therefore, the camera location is $O = -R^T T$, and the corresponding view ray is along $\overrightarrow{OW_c}$. The rotation angle $r(x_c, y_c)$ is the angle between the unit vector $(0, 1)$ along $y$ direction and $\overrightarrow{OW_c}$. After the rotation angle map $r(x_c, y_c)$ is calculated, it is then quantized by $q$ degree into the rotation range $\{r_0, r_1, ..., r_m\}$. The rotation selection mask for rotation angle $r_i$ is $\mathbb{1}(r = r_i)$, and the fused features are $F = \sum_i \mathbb{1}(r = r_i) \otimes F_i$. Examples of the rotation selection masks are presented in the second row of Fig. 12.

| Dataset | resolution | view | train / test | crowd |
|---------|-----------|------|-------------|-------|
| PETS2009 | 768×576 | 3 | 1105 / 794 | 20-40 |
| DukeMTMC | 1920×1080 | 4 | 700 / 289 | 10-30 |
| CityStreet | 2704×1520 | 3 | 300 / 200 | 70-150 |

**Table 2** The comparison of three multi-view datasets.

4.3 Training details

A two-stage procedure is used to train the model. The first stage trains the main scene-level density map estimation task as well as auxiliary view-level density map estimation tasks. The auxiliary task for late fusion consists of auxiliary losses applied between the predicted and GT density maps for each view. The auxiliary task for early fusion uses an auxiliary branch of 3 layers of FCN to predict density maps for each view, followed by an auxiliary loss on the view-level predicted density maps. The learning rate is set to 1e-4. In the second stage, the auxiliary view-level density map estimation tasks are removed, leaving only the scene-level task. FCN-7 (either density map estimator or feature extractor) is fixed and the fusion and scale selection parts are trained. The loss function is the pixel-wise squared error between the ground-truth and predicted density maps. The learning rate is set to 1e-4, and decreases to 5e-5 during training. After training the two stages, the model is fine-tuned end-to-end. The training batch-size is set to 1 in all experiments.

# 5 Datasets and Experiment Setup

In this section we introduce the 3 multi-view counting datasets and the experiment settings.

5.1 Datasets

We test the proposed multi-view counting framework on two existing datasets, PETS2009 and DukeMTMC, and our newly collected CityStreet dataset. Table 2 provides a summary, and Fig. 13 shows examples.

5.1.1 PETS2009

PETS2009 (Ferryman and Shahrokni 2009) is a multi-view sequence dataset containing crowd activities from 8 views. The first 3 views are used for the experiments, as the other 5 views have low camera angle, poor image quality, or unstable frame rate. To balance the crowd levels, we use sequences S1L3 (14_17, 14_33), S2L2 (14_55) and S2L3 (14_41) for training (1105 images in total), and S1L1 (13_57, 13_59), S1L2 (14_06, 14_31) for testing (794 images). The calibration parameters (extrinsic and intrinsic) for the cameras are provided with the dataset.

To obtain the annotations across all views, we use the View 1 annotations provided by Leal-Taixé et al. (2015) and project them to other views, followed by manual annotations to get all the people heads in the images.

5.1.2 DukeMTMC

DukeMTMC (Ristani et al. 2016) is a multi-view video dataset for multi-view tracking, human detection, and ReID. The multi-view video dataset has videos from 8 synchronized cameras for 85 minutes with 1080p resolution at 60 fps. For our counting experiments, we use 4 cameras that have overlapping fields-of-view (cameras 2, 3, 5 and 8). The synchronized videos are sampled every 3 seconds, resulting in 989 multi-view images. The first 700 images are used for training and the remaining 289 for testing. Using the same sampling method for creating the DukeMTMC training set, 200 multi-view frames are extracted from DukeMTMC Test Hard set for extra evaluation. Camera extrinsic and homography parameters are provided by the dataset. In the original dataset, annotations for each view are only provided in the view ROIs, which are all non-overlapping on the ground-plane and in Camera 8, only R3 region is used and R1 and R2 are excluded. Since we are interested in overlapping cameras, we project the annotations from each camera view to the overlapping areas in all other views. Region R2 (see Fig. 13) is excluded during the experiment, since there are no annotations provided there.

5.1.3 CityStreet

We collected a multi-view video dataset of a busy city street in Hong Kong using 5 synchronized cameras. The videos are about 1 hour long with 2.7k (2704×1520) resolution at 30 fps. We select Cameras 1, 3 and 4 for the experiment (see Fig. 13 bottom). The cameras' intrinsic and extrinsic parameters are estimated using the calibration algorithm from Zhang (2000). 500 multi-view images are uniformly sampled from the videos, and the first 300 are used for training and remaining 200 for testing. The ground-truth 2D and 3D annotations are obtained as follows. The head positions of the first camera-view are annotated manually, and then projected to other views and adjusted manually. Next, for the second camera view, new people (not seen in the first view), are also annotated and then projected to the other views. This process is repeated until all people in the scene are annotated and associated across all camera views. Our dataset has larger crowd numbers (70-150), compared with PETS (20-40) and Duke-MTMC (10-30). Our dataset also contains more crowd

**Fig. 13** Examples from 3 multi-view counting datasets. The first column shows the camera frames and annotations. The second column shows the camera layout and scene-level ground-plane density maps. Note that 'R2' region of Camera 8 of DukeMTMC dataset is not used since no annotations are available.

scale variations and occlusions due to vehicles and fixed structures.

## 5.2 Experiment setup

### 5.2.1 Ground-truth settings

The ground-truth scene-level density maps are created by convolving the people's ground-plane annotation map with a fixed-bandwidth Gaussian kernel. The people's ground-plane annotations are estimated from their camera-view annotations and the camera calibration information. First, each person's height in 3D is estimated by selecting a height, among a candidate set in the range [1.6m, 2.0m] (step of 0.1m), so as to minimize the differences among the ground-plane coordinates obtained by projecting that person's 2D camera-view annotations to the 3D world. Second, with the estimated person's height, a person's ground-plane annotation is set to the average of that person's 2D annotations projected to the 3D world. To construct the ground-truth scene-level density maps, we follow single-image counting and choose a fixed-bandwidth $\sigma = 3$ for the Gaussian kernel.

The image resolutions ($w \times h$) used in the experiments are: 384×288 for PETS2009, 640×360 for Duke-MTMC, and 676×380 for CityStreet. The resolutions of the scene-level ground-plane density maps are: 152×177 for PETS2009, 160×120 for DukeMTMC and 160×192 for CityStreet. For the detection baseline, the original

image resolutions are used (Faster-RCNN will resize the images).

### 5.2.2 Methods

We test our multi-view fusion models, denoted as "Late fusion", "Naïve early fusion", "MVMS" (multi-view multi-scale early fusion), and "MVMSR" (MVMS with rotation selection). The late fusion model uses projection normalization. MVMS uses learnable scale selection, and a 3-scale image pyramid with zoom factor of 0.5. MVMSR uses 3 rotation selection layers with filter number $F = 32$ and quantization angle $Q = 10°$, $45°$ and $45°$ for PETS2009, DukeMTMC and CityStreet, respectively. Besides, we have also performed the models with different feature extraction backbones. These settings will be tested later in the ablation study.

For comparisons, we test and compare with several comparison methods. The first comparison method is an approach to fusing camera-view density maps into a scene-level crowd count, denoted as "Dmap weighted", which is an adaptation from Ryan et al. (2014). First single image counting model is applied to get the density map $D_i$ for each camera-view. The density maps are then fused into a scene-level count using a weight map $W_i$ for each view,

$$C = \sum_i \sum_{x_0, y_0} W_i(x_0, y_0) D_i(x_0, y_0), \qquad (9)$$

where the summations are over the camera-views and the image pixels. The weight map $W_i$ is constructed

| Dataset | Method | MSE | NAE | MAE/GAME(0) | GAME(1) | GAME(2) |
|---|---|---|---|---|---|---|
| PETS2009 | Dmap weighted | 7.29 | 0.182 | 5.62 | - | - |
| | Detection+ReID | 7.01 | 0.174 | 5.46 | - | - |
| | Feature concatenation | 8.96 | 0.300 | 7.32 | - | - |
| | Stitching | 14.51 | 0.337 | 10.90 | - | - |
| | Late fusion | 5.18 | 0.138 | 3.92 | 6.38 | 7.75 |
| | Naïve early fusion | 6.76 | 0.199 | 5.42 | 7.26 | 8.13 |
| | MVMS | **4.83** | **0.124** | **3.49** | **5.30** | **6.27** |
| | MVMSR | 4.93 | 0.130 | 3.62 | 5.37 | 6.98 |
| DukeMTMC | Dmap weighted | 2.06 | 0.186 | 1.54 | - | - |
| | Detection+ReID | 2.30 | 0.355 | 1.89 | - | - |
| | Feature concatenation | 5.23 | 1.030 | 4.44 | - | - |
| | Stitching | 1.65 | 0.215 | 1.24 | - | - |
| | Late fusion | 1.63 | 0.187 | 1.29 | 1.77 | 2.22 |
| | Naïve early fusion | 1.90 | 0.199 | 1.47 | 2.00 | 2.62 |
| | MVMS | 1.28 | 0.122 | 0.95 | 1.24 | 1.50 |
| | MVMSR | **1.17** | **0.118** | **0.89** | **1.19** | **1.42** |
| CityStreet | Dmap weighted | 11.46 | 0.120 | 9.36 | - | - |
| | Detection+ReID | 21.18 | 0.193 | 17.48 | - | - |
| | Feature concatenation | 21.34 | 0.245 | 18.33 | - | - |
| | Stitching | 10.55 | 0.107 | 8.76 | - | - |
| | Late fusion | 9.63 | 0.099 | 8.06 | 12.75 | 23.10 |
| | Naïve early fusion | 9.85 | 0.100 | 8.11 | 12.73 | 22.93 |
| | MVMS | 9.02 | 0.096 | 7.36 | 11.95 | 20.44 |
| | MVMSR | **8.49** | **0.086** | **6.98** | **11.39** | **19.79** |

**Table 3** The scene-level counting performance of different methods on the 3 datasets. MSE, NAE, MAE and GAME are used as evaluation metrics. Note that for comparison methods whose outputs are not density maps, GAME(1) and GAME(2) are not applicable. FCN-7 is used as feature backbone for PETS2009, and CSR-net is used as feature backbone for CityStreet and DukeMTMC. For MVMSR, the filter number is 32, the layer number is 3 and the quantization angle is 10°, 45° and 45° for the PETS2009, DukeMTMC and CityStreet, respectively.

based on how many views can see a particular pixel. In other words, $W_i(x_0, y_0) = 1/t$, where $t$ is the number of views that can see the projected point $\mathcal{P}(x_0, y_0)$. Note that Ryan et al. (2014) used this simple fusion approach with traditional regression-based counting (in their setting, the $D_i$ map is based on the predicted counts for crowd blobs). We also test on the methods with different single-image counting models. Here, we are using recent DNN-based methods (CSR-net) and crowd density maps, which outperform traditional regression-based counting, and hence form a stronger baseline method compared to Ryan et al. (2014).

The second comparison method is using human detection methods and person re-identification (ReID), denoted as "Detection + ReID". First, Faster-RCNN (Ren et al. 2015) is used to detect humans in each camera-view. Next, the scene geometry constraints and the ReID method Circle2020 ((Sun et al. 2020; Wang et al. 2018)) are used to associate the same people across views. Specifically, each detection box's top-center point in one view is projected to other views, and ReID is performed between the original detection box and detection boxes near the projected point in other views. Finally, the scene-level people count is obtained by counting the number of unique people among the detection boxes in all views. The bounding boxes needed for training are created with the head annotations and the perspective map of each view.

The third comparison method is to simply concatenate the features of the different views and directly regress a scene-level count. We've used the LCC (Liu et al. 2020) as the feature extractor, where the output of the scale-aware module of LCC is used as extracted features.

The fourth comparison method is to "stitch" together the counts from different camera views. In particular, first a set of non-overlapping ROIs on the ground-plane are formed by assigning ground-plane pixels to the closest camera. These ROIs are then projected into their camera views. Next, single-view counting (CSR-net) is performed on each camera-view, and the ROI count in each camera view is obtained. Finally, the scene-level count is the sum of the ROI counts.

*5.2.3 Evaluation*

The mean absolute error (MAE), mean squared error (MSE), and normalized (relative mean) absolute error (NAE) are used to evaluate multi-view counting performance, comparing the scene-level predicted counts and the ground-truth scene-level counts. Besides, Grid Average Mean absolute Error (GAME) (Guerrero-Gómez-Olmedo et al. 2015) is also used to evaluate the local

counting performance of the predicted scene-level density maps of the proposed methods. The definitions of these evaluation metrics are as follows.

$$MAE = \frac{1}{N}\sum_i^N |c_i - \hat{c}_i|, \tag{10}$$

$$MSE = \sqrt{\frac{1}{N}\sum_i^N (c_i - \hat{c}_i)^2}, \tag{11}$$

$$NAE = \frac{1}{N}\sum_i^N |c_i - \hat{c}_i|/\hat{c}_i, \tag{12}$$

$$GAME(L) = \frac{1}{N}\sum_i^N (\sum_{l=1}^{4^L} |c_i^l - \hat{c}_i^l|), \tag{13}$$

where $N$ is the number of the test images, $c_i$ and $\hat{c}_i$ are the estimated and ground-truth people count in the $i$-th image. As to GAME metric, the scene-level density maps are divided in $4^L$ non-overlapping patches and compute the average of the $MAE$ of these patches. $c_i^l$ and $\hat{c}_i^l$ are the estimated and ground-truth people count of the patch $l$ of $i$-th image. Note that $MAE$ equals the $GAME$ when $L = 0$.

In addition, we also evaluate the predicted counts in each camera-view. The ground-truth count for each camera-view is obtained by summing the ground-truth scene-level density map over the region covered by the camera's field-of-view. Note that people that are totally occluded from the camera, but still within its field-of-view, are still counted.

## 6 Experiment Results

In this section, the scene-level counting performance of the proposed DNN-based multi-view fusion methods are evaluated against other multi-camera counting methods. We also demonstrate the single-view counting performance using multi-view cameras. In terms of both evaluation perspectives, the proposed method can achieve better counting results on all 3 multi-view counting datasets.

### 6.1 Scene-level counting performance

In this section, we test the proposed multi-view counting models in terms of scene-level counting performance on the 3 multi-view counting datasets, PETS2009 (Ferryman and Shahrokni 2009), DukeMTMC (Ristani et al. 2016) and CityStreet. The results are presented in Tables 3, 7, 8, and 9, and examples shown in Fig. 14.

| Method | MSE | NAE | MAE |
|---|---|---|---|
| Hybrid (Dittrich et al. 2017) | - | - | 2.03 |
| Late fusion (w/ PN) | 2.56 | 0.241 | 1.53 |
| Naive early fusion | 2.14 | 0.203 | 1.71 |
| MVMS | 1.26 | 0.091 | 0.98 |
| MVMSR | **1.24** | **0.089** | **0.96** |

**Table 4** Extra experiment results on PETS S1L1 (views 1 and 2) comparing with a traditional multi-view method. 'PN' means projection normalization.

| Method | MSE | NAE | MAE |
|---|---|---|---|
| Dmap weighted | 7.21 | 0.437 | 4.17 |
| Detection+ReID | 7.06 | 0.371 | 3.71 |
| Late fusion (w/ PN) | **4.82** | **0.307** | **2.62** |
| Naïve early fusion | 6.13 | 0.361 | 2.82 |
| MVMS | 6.20 | 0.328 | 2.81 |
| MVMSR | 6.00 | 0.329 | 2.89 |

**Table 5** Extra experiment results on DukeMTMC Test Hard set.

*6.1.1 PETS2009*

The scene-level counting results on PETS2009 are shown in Table 3 (top row) and Table 7 ("Scene" column). On PETS2009, our proposed multi-view fusion models (use FCN-7 as backbone) achieve better results than the two comparison methods. Detection+ReID (Circle2020) performs worst on this dataset because the people are close together in a crowd, and occlusion causes severe misdetection. Among our three multi-view fusion models, naïve early fusion performs worse, which suggests that the scale variations in multi-view images limits the performance. Furthermore, MVMS performs much better than other models, which shows the multi-scale framework with scale selection strategies can improve the feature-level fusion to achieve better performance.

The performance of using one camera for the scene-level counting task is not as good as using multi-cameras. In particular, using Dmap with cameras 1 or 2 performs poorly due to the limited field-of-view. Dmap using camera 3 achieves slightly better performance than using multi-cameras (Dmap weighted, CSR-net backboned) because most people can already be seen in camera 3. This also suggests that "Dmap weighted" cannot fuse the multi-view information well. Nonetheless, our fusion methods all outperform the Dmap weighted, demonstrating the efficacy of the projection and fusion stages. The feature concatenation method outputs a scene-level count, and does not consider the geometry relationship of the camera views, thus achieving worse results than the proposed methods. The Stitching method, which only considers the distance between objects and cameras but neglects the occlusions in the camera views, also performs worse than our methods. Finally, using multiple cameras to count with "Detec-

tion+ReID" improves the performance over single cameras, but still has higher error than our fusion methods.

**PETS S1L1.** We next compare our method with a traditional multi-view counting method "Hybrid" (Dittrich et al. 2017) on the subset of PETS2009 dataset, PETS S1L1, which is presented in Table 4. Dittrich et al. (2017) proposed two approaches (head detector and count regression) by fusing hand-crafted features (corner points or Harr feature) from multiple cameras for multi-view counting. Similar to Dittrich et al. (2017), we use PETS2009 S1L1 13_57 (view 1 and 2) for training and 13_59 (view 1 and 2) for testing. Our fusion models (FCN-7 backbone) all achieve better performance than the multi-view counting method based on traditional hand-crafted low-level features, and MVMSR achieves the best scene-level counting performance.

### 6.1.2 DukeMTMC

The scene-level counting results on DukeMTMC are shown in Table 3 (middle row) and Table 8 ("Scene" column). On DukeMTMC, our multi-view fusion models (use CSR-net as backbone) achieve better performance than comparison methods at the scene-level counting task. Due to lower crowd numbers in DukeMTMC, the performance gap among the 3 fusion models is not large – but MVMS and MVMSR still perform best and MVMSR is better than MVMS. Furthermore, results from comparison methods also show that using a single camera is not adequate for the scene-level counting task. Dmap with a single camera 8 performs slightly better than using multi-cameras (Dmap weighted) due to the large field-of-view of camera 8, and the limitations of the weighted fusion. Since camera views in DukeMTMC dataset share smaller area of overlapping regions and the occlusion issue is not severe, the Stitching method performs relatively better than other comparison methods, but the proposed MVMS and MVMSR still perform better.

**DukeMTMC Test Hard.** Finally, the scene-level counting results on the DukeMTMC Test Hard set, which contains more crowds, are presented in Table 5. Our fusion model (FCN-7 backbone) achieves better scene-level counting results than the baselines. Among our methods, late fusion has slightly lower error than MVMS/MVMSR.

### 6.1.3 CityStreet

The scene-level counting results on CityStreet are shown in 3 (bottom row) Table 9 ("Scene" column). On CityStreet, our multi-view fusion models achieve better results than the comparison methods. Compared to PETS2009,

| Dataset | Method | MSE | NAE | MAE |
|---------|--------|-----|-----|-----|
| PETS2009 | CVF (Zheng et al. 2021) | - | - | **3.08** |
| | CVCS (Zhang et al. 2021) | - | 0.165 | 5.17 |
| | MVMS | **4.83** | **0.124** | 3.49 |
| | MVMSR | 4.93 | 0.130 | 3.62 |
| DukeMTMC | CVF (Zheng et al. 2021) | - | - | **0.87** |
| | CVCS (Zhang et al. 2021) | - | 0.525 | 2.83 |
| | MVMS | 1.28 | 0.122 | 0.95 |
| | MVMSR | **1.17** | **0.118** | 0.89 |
| CityStreet | CVF (Zheng et al. 2021) | - | - | 7.08 |
| | CVCS (Zhang et al. 2021) | - | 0.117 | 9.58 |
| | MVMS | 9.02 | 0.096 | 7.36 |
| | MVMSR | **8.49** | **0.086** | **6.98** |

**Table 6** The scene-level counting performance of different methods on the 3 datasets. MSE, NAE and MAE are used as evaluation metrics. Note that only MAE is provided in CVF (Zheng et al. 2021), and MAE and NAE are provided in CVCS (Zhang et al. 2021).

CityStreet has larger crowds and more occlusions and scale variations. Therefore, the performances of the baseline methods decreases significantly, especially Detection+ReID. Due to large camera angle change and severe occlusions in the CityStreet dataset, Feature concatenation and Stitching cannot perform well on the larger and more complicated dataset. Our MVMSR model achieves much better performance on CityStreet than all other models. The reason is the 3 views of the CityStreet dataset have larger view angle change than the other two datasets, which can better demonstrate the effectiveness of the rotation selection in the multi-view fusion process. Furthermore, similar to the other two datasets, using multi-cameras achieves better scene-level counting performance than using a single camera.

### 6.1.4 Comparison with concurrent methods

We next compare with two recent multi-view counting methods published concurrently during the revision of our paper: CVF (Zheng et al. 2021) and CVCS (Zhang et al. 2021). The comparison is shown in Table 6, and our proposed method achieves the best performance on the largest dataset CityStreet.

### 6.2 Single-view counting performance

We next evaluate the single-view counting performance, which is the people count within the single-camera's field-of-view. Here we mainly aim to show that multi-view information can improve single-view counting over using a single camera. Note, the comparison method feature concatenate and stitching directly output a scene-level count and are not used for comparison in this section.

*PETS2009:* The single-view counting results on PETS2009 are shown in Table 7. Columns 'C1', 'C2', and 'C3' correspond for single-view counting in regions within the

| Method (camera) | PETS2009 (Ferryman and Shahrokni 2009) | | | |
| | Scene | C1 | C2 | C3 |
|---|---|---|---|---|
| Dmap (camera 1) | 13.74/0.413/12.19 | 4.55/0.213/3.96 | - | - |
| Dmap (camera 2) | 13.48/0.404/12.39 | - | 9.43/0.309/8.33 | - |
| Dmap (camera 3) | 7.95/0.239/6.89 | - | - | 6.37/0.201/5.46 |
| Dmap weighted (multiview) | 7.29/0.182/5.62 | 4.02/0.169/3.61 | 4.25/0.136/3.42 | 5.21/0.149/4.23 |
| Detection+ReID (camera 1) | 17.99/0.545/17.24 | 9.75/0.356/8.57 | - | - |
| Detection+ReID (camera 2) | 16.27/0.485/15.54 | - | 12.33/0.393/11.19 | - |
| Detection+ReID (camera 3) | 16.27/0.503/16.33 | - | - | 15.60/0.472/14.59 |
| Detection+ReID (multiview) | 7.01/0.174/5.46 | 8.22/0.238/6.55 | 9.33/0.300/7.09 | 13.50/0.400/11.76 |
| Late fusion (multiview) | 5.18/0.138/3.92 | 3.20/0.143/2.62 | 4.19/0.137/3.17 | 5.00/0.150/3.97 |
| Naïve (multiview) | 6.76/0.199/5.42 | 3.13/0.124/2.37 | 5.76/0.179/4.27 | 6.36/0.192/4.92 |
| MVMS (multiview) | **4.83/0.124/3.49** | 2.22/0.084/1.66 | 3.67/0.103/2.58 | **4.58/0.127/3.46** |
| MVMSR (multiview) | 4.93/0.130/3.62 | **2.17/0.077/1.57** | **3.30/0.097/2.38** | 4.76/0.133/3.64 |

**Table 7** Comparison of the scene-level (left) and the single-view counting (right) measured with mean square error, mean absolute error and relative mean absolute error (MSE/NAE/MAE) on PETS2009. Column "Scene" denotes the scene-level counting error. Columns "C1", "C2" and "C3" refer to the single-view counting error for the region within the field-of-view of cameras 1, 2 and 3. "camera" indicates the camera(s) used for counting. The late fusion model uses projection normalization, and MVMS and MVMSR uses learnable scale selection and FCN-7 is used as the feature backbone.

| Method (camera) | DukeMTMC (Ristani et al. 2016) | | | | |
| | Scene | C2 | C3 | C5 | C8 |
|---|---|---|---|---|---|
| Dmap (camera 2) | 6.07/0.613/5.19 | 0.97/0.487/0.73 | - | - | - |
| Dmap (camera 3) | 8.73/1.000/8.03 | - | 1.28/0.647/0.79 | - | - |
| Dmap (camera 5) | 7.39/0.830/6.72 | - | - | 0.81/0.575/0.49 | - |
| Dmap (camera 8) | 2.43/0.258/1.87 | - | - | - | 1.57/0.232/1.21 |
| Dmap weighted | 2.71/0.250/2.11 | 1.35/0.426/1.02 | 1.28/0.647/0.79 | 1.42/0.663/0.89 | 1.83/0.201/1.30 |
| Det+ReID (camera 2) | 4.42/0.425/3.51 | 2.28/1.03/2.06 | - | - | - |
| Det+ReID (camera 3) | 7.93/0.890/7.20 | - | 0.55/**0.132**/0.25 | - | - |
| Det+ReID (camera 5) | 7.11/0.782/6.38 | - | - | 1.29/0.524/0.96 | - |
| Det+ReID (camera 8) | 5.85/0.620/5.10 | - | - | - | 4.28/0.541/3.58 |
| Det+ReID (multiview) | 2.30/0.355/1.89 | 0.94/0.513/0.75 | 0.71/0.584/0.40 | 2.37/1.09/1.89 | 3.62/0.422/2.86 |
| Late fusion (multiview) | 1.63/0.187/1.29 | 0.64/0.263/0.45 | 0.56/1.040/0.35 | 0.66/0.548/0.37 | 1.48/0.198/1.15 |
| Naïve (multiview) | 1.90/0.199/1.47 | 0.59/0.265/0.44 | 0.66/0.970/0.44 | 0.91/0.768/0.60 | 1.64/0.195/1.23 |
| MVMS (multiview) | 1.28/0.122/0.95 | 0.50/0.199/0.33 | 0.40/0.677/0.21 | 0.59/**0.401**/0.31 | 1.11/0.125/0.81 |
| MVMSR (multiview) | **1.17/0.118/0.89** | **0.45/0.183/0.30** | **0.38**/0.758/**0.19** | **0.54**/0.410/**0.29** | **1.02/0.118/0.76** |

**Table 8** Comparison of the scene-level (left) and the single-view counting (right) using mean square error, mean absolute error and relative mean absolute error (MSE/NAE/MAE) on DukeMTMC. CSR-net is used as the feature backbone. See the caption of Table 7 for further description.

fields-of-view of cameras 1, 2, and 3, respectively. On PETS2009, our 3 multi-view fusion models can achieve better results than the two comparison methods in terms of all single-camera counting, which demonstrates that the proposed multi-view fusion DNNs can well integrate the information from multi-views to improve the counting performance in different regions. Furthermore, MVMS performs much better than other models, which also shows the multi-scale framework with scale selection strategies can improve the feature-level fusion to achieve better performance. Finally, the comparison methods' single-view counting performance can also be improved with the aid of other cameras.

*DukeMTMC:* On DukeMTMC, our multi-view fusion models can achieve better performance than comparison methods on most camera-views (see Table 8). Detection+ReID achieves the good result on camera 3 because this camera is almost parallel to the horizontal plane, has low people count, and rarely has occlusions, which is an ideal operating regime for the detector. Fi-

nally, the single-view counting performance is mostly improved with the aid of multi-cameras.

*CityStreet:* On CityStreet (see Table 9), our 3 multi-view fusion models achieve better results than the comparison methods. Due to severe occlusions and scale changes, the Detection+ReID methods perform badly, even with the aid of other cameras. However, using multi-cameras still improves the single-view counting performance, and our methods are the most effective at multi-view fusion.

### 6.3 Ablation studies

We next present ablation studies on the various components of our fusion pipeline.

### 6.3.1 Backbone for MVMS and MVMSR

We compare different feature extraction backbones in the proposed multi-view counting framework: FCN-7,

| Dataset | CityStreet | | | |
|---|---|---|---|---|
| Method (camera) | Scene | C1 | C3 | C4 |
| Dmap (camera 1) | 11.70/0.110/9.31 | 10.21/0.112/8.51 | - | - |
| Dmap (camera 3) | 11.64/0.199/9.41 | - | 11.83/0.129/9.23 | - |
| Dmap (camera 4) | 24.24/0.256/21.92 | - | - | 22.84/0.240/20.30 |
| Dmap weighted (multiview) | 11.46/0.120/9.36 | 9.30/0.101/7.87 | 11.19/0.121/9.19 | 12.84/0.116/10.16 |
| Detection+ReID (camera 1) | 49.68/0.513/45.80 | 45.71/0.483/41.38 | - | - |
| Detection+ReID (camera 3) | 45.09/0.453/40.87 | - | 37.94/0.391/32.94 | - |
| Detection+ReID (camera 4) | 35.10/0.323/30.03 | - | - | 33.16/0.311/28.57 |
| Detection+ReID (multiview) | 21.18/0.193/17.48 | 42.36/0.456/37.76 | 35.10/0.355/29.27 | 22.84/0.228/18.21 |
| Late fusion | 9.63/0.099/8.06 | 9.71/0.110/8.36 | 9.51/0.110/7.99 | 9.01/0.089/7.46 |
| Naïve early fusion | 9.85/0.100/8.11 | 10.04/0.108/8.35 | 9.42/0.106/7.74 | 9.65/0.098/7.94 |
| MVMS (multiview) | 9.02/0.096/7.36 | 9.59/0.110/7.87 | 8.44/0.100/6.87 | **7.59**/0.081/**6.24** |
| MVMSR (multiview) | **8.49/0.086/6.98** | **8.51/0.094/7.05** | **8.06/0.089/6.49** | 7.89/**0.078**/6.44 |

**Table 9** Comparison of the scene-level and the single-view counting using mean square error, mean absolute error and relative mean absolute error (MSE/NAE/MAE) on CityStreet. CSR-net is used as the feature backbone. See the caption of Table 7 for further description.

| Backbone | Method | MSE | NAE | MAE/GAME(0) | GAME(1) | GAME(2) |
|---|---|---|---|---|---|---|
| FCN-7 | Late fusion | 10.24 | 0.097 | 8.12 | **13.05** | 21.14 |
| | Naïve early fusion | 10.11 | 0.096 | 8.10 | 13.06 | **20.98** |
| | MVMS | 10.05 | 0.096 | 8.01 | 13.67 | 21.99 |
| | MVMSR | **9.73** | **0.090** | **7.37** | 13.89 | 22.60 |
| CSR-net | Late fusion | 9.63 | 0.099 | 8.06 | 12.75 | 23.10 |
| | Naïve early fusion | 9.85 | 0.100 | 8.11 | 12.73 | 22.93 |
| | MVMS | 9.02 | 0.096 | 7.36 | 11.95 | 20.44 |
| | MVMSR | **8.49** | **0.086** | **6.98** | **11.39** | **19.79** |
| LCC | Late fusion | 10.51 | 0.113 | 8.71 | 15.45 | 26.29 |
| | Naïve early fusion | 9.96 | 0.103 | 7.97 | 12.99 | 22.56 |
| | MVMS | 9.86 | 0.093 | 7.67 | 13.92 | 22.97 |
| | MVMSR | **9.46** | **0.086** | **7.42** | **12.77** | **21.22** |

**Table 10** Comparison of different backbones for scene-level counting performance on CityStreet, where the settings for the rotation module are the same: the filter number is 32, the layer number is 3 and the quantization angle is 45°.

CSR-net (Li et al. 2018) (first 7 layers of VGG) and LCC (Liu et al. 2020) (use the output of the scale-aware module as extracted features).

First, the counting results of the 4 fusion models on CityStreet are presented in Table 10. Generally, using larger backbone performs better than FCN-7. However, the performance gap using different is larger for Dmap weighted, indicating that the single-view density maps of CSR-Net are more accurate. This suggests that the scale-selection module in MVMS is sufficient for handling scale changes in multi-view counting, and the benefits of using dilated convolutions to handle single-view scale changes are diminished. Furthermore, the improvement of MVMSR and MVMS over late/naïve fusion is consistent among the 3 backbones. Finally, the proposed MVMSR method achieves the best performance on all 3 backbones, which indicates the effectiveness of the rotation selection module.

Second, the counting results of MVMS and MVMSR on the 3 datasets are presented in Table 11. Generally, with larger backbone, the performance of MVMS or MVMSR can be improved on the CityStreet and Duke-MTMC dataset, and MVMSR is better than MVMS. On PETS2009 dataset, the smaller backbone performs

better than larger backbones, and the possible reason is the larger backbone models are overfitting on the dataset since more CNNs layers are used. Furthermore, on PETS2009, since the camera angle change is not large enough and most people can be seen by camera 3, the best performance of MVMSR is slightly worse than MVMS. Nonetheless, MVMSR is better than MVMS on the more complicated datasets, CityStreet and Duke-MTMC.

### 6.3.2 Normalization in the late fusion model

We perform an ablation study on the late fusion model (FCN-7 backbone) with and without the projection normalization step, and the results are presented in Table 12 (top). Using projection normalization reduces the error of the late fusion model, compared to not using the normalization step. This demonstrates the importance of maintaining the total count when projecting the density map into the ground-plane representation.

### 6.3.3 Scale selection in MVMS

We perform an ablation study on the scale-selection strategy of MVMS, and the results are presented in Ta-

| Backbone | Method | MSE | NAE | MAE/GAME(0) | GAME(1) | GAME(2) |
|---|---|---|---|---|---|---|
| | MVMS(FCN-7) | **4.83** | **0.124** | **3.49** | **5.30** | **6.27** |
| | MVMS(CSR-net) | 5.26 | 0.140 | 3.99 | 6.27 | 7.51 |
| PETS2009 | MVMS (LCC) | 5.62 | 0.151 | 4.36 | 6.68 | 7.98 |
| | MVMSR (FCN-7) | 4.93 | 0.130 | 3.62 | 5.37 | 6.98 |
| | MVMSR (CSR-net) | 5.51 | 0.140 | 4.15 | 6.56 | 8.30 |
| | MVMSR (LCC) | 5.60 | 0.151 | 4.37 | 5.80 | 6.94 |
| | MVMS(FCN-7) | 1.24 | 0.170 | 1.03 | 1.53 | 1.92 |
| | MVMS(CSR-net) | 1.28 | 0.122 | 0.95 | 1.24 | 1.50 |
| DukeMTMC | MVMS (LCC) | 1.38 | 0.132 | 1.04 | 1.26 | 1.49 |
| | MVMSR (FCN-7) | 1.31 | 0.144 | 1.01 | 1.50 | 2.02 |
| | MVMSR (CSR-net) | **1.17** | **0.118** | **0.89** | **1.19** | **1.42** |
| | MVMSR (LCC) | 1.26 | 0.129 | 0.94 | 1.29 | 1.57 |
| | MVMS(FCN-7) | 10.05 | 0.096 | 8.01 | 13.67 | 21.99 |
| | MVMS(CSR-net) | 9.02 | 0.096 | 7.36 | 11.95 | 20.44 |
| CityStreet | MVMS (LCC) | 9.86 | 0.093 | 7.67 | 13.92 | 22.97 |
| | MVMSR (FCN-7) | 9.73 | 0.090 | 7.37 | 13.89 | 22.60 |
| | MVMSR (CSR-net) | **8.49** | **0.086** | **6.98** | **11.39** | **19.79** |
| | MVMSR (LCC) | 9.46 | 0.086 | 7.42 | 12.77 | 21.22 |

**Table 11** Comparison of different backbones for MVMS/MVMSR on PETS2009, DukeMTMC and CityStreet. For MVMSR, the filter number is 32, the layer number is 3 and the quantization angle is 10°, 45° and 45° for PETS2009, DukeMTMC and CityStreet, respectively.

| Dataset | PETS2009 | | | | DukeMTMC | | | | | CityStreet | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Region | C1 | C2 | C3 | Scene | C2 | C3 | C5 | C8 | Scene | C1 | C3 | C4 | Scene |
| Late fusion (w/ PN) | **2.62** | **3.17** | **3.97** | **3.92** | 0.49 | 0.77 | **0.39** | 1.15 | 1.27 | 8.14 | 7.72 | 8.08 | 8.12 |
| Late fusion (w/o PN) | 2.75 | 3.86 | 4.37 | 4.22 | 0.63 | **0.73** | 0.51 | 1.31 | 1.43 | 9.89 | 9.60 | 9.82 | 9.87 |
| MVMS (fixed) | 1.74 | **2.57** | 3.81 | 3.82 | 0.65 | **0.46** | **0.88** | 1.44 | 1.09 | 8.11 | 7.83 | 8.32 | **7.80** |
| MVMS (learnable) | **1.66** | 2.58 | **3.46** | **3.49** | **0.63** | 0.52 | 0.94 | **1.36** | **1.03** | **7.99** | **7.63** | **7.91** | 8.01 |

**Table 12** Ablation study (MAE) comparing the late fusion model with and without projection normalization (PN), and MVMS with fixed or learnable scale selection.

| Dataset | fixed discrete | fixed soft | learnable soft |
|---|---|---|---|
| PETS2009 | 3.82 | 3.59 | 3.49 |
| CityStreet | 7.80 | 8.55 | 8.01 |

**Table 13** MAE comparison of MVMS model (FCN-7 backbone) selection module settings: fixed scale selection with discrete mask or soft-mask, and learnable scale selection with soft-mask.

ble 12 (bottom). Most of the time the learnable scale-selection strategy can achieve lower error than fixed scale-selection. We note that using MVMS with fixed scale-selection strategy still outperforms the naïve early fusion, which performs no scale selection. Thus obtaining features that have consistent scales across views is an important step when fusing the multi-view feature maps.

The proposed two scale selection modules use different mask methods: discrete mask for fixed scale selection and soft-mask for learnable scale selection. We perform another ablation study on fixed scale selection using the soft-mask, which is presented in Table 13. When using soft-masks, learnable scale selection still outperforms fixed scale selection.

*6.3.4 Rotation module in MVMSR*

We perform the ablation study on the rotation module of MVMSR on the CityStreet dataset (with CSR-net

backbone), including the number of filters $F$, the number of layers $L$, and the rotation quantization angle $Q$ of the rotation selection layer.

**Number of filters** $F$**:** The results of the ablation study on the number of filters $F$ is shown in Table 14, where the number of layers in rotation selection is 3 and the rotation quantization angle is 45°. Increasing the number of filters does not necessarily improve the performance of the scene-level counting performance. The reason is because the rotation selection layer naturally handles the effect of feature rotations in the projection step, and thus less filters are required when compared to without the rotation selection, where each rotation needs a separate filter. Besides, using more filters decreases the speed of the model in the inference stage. Therefore, we choose $F = 32$ in the remaining experiments, whose performance is better compared to others.

**Number of layers** $L$**:** We perform an ablation study on the number of layers $L$ of the rotation module in Table 15, where the number of filters is 32 and the rotation quantization angle is 45°. From the table, we conclude that fewer rotation selection layers may not reduce the feature rotation effect and too many layers may make the model overfit and decrease the counting performance. The choice of $L = 5$ achieves the best

| $F$ | MAE | MSE | NAE | G(1) | G(2) | FPS |
|---|---|---|---|---|---|---|
| 8 | 7.56 | 9.36 | 0.091 | 11.54 | 20.70 | **5.8** |
| 16 | 7.24 | 8.74 | 0.088 | 11.73 | 19.81 | 5.6 |
| 32 | **6.98** | **8.49** | **0.086** | **11.39** | **19.79** | 5.5 |
| 64 | 7.09 | 8.73 | 0.088 | 11.90 | 20.47 | 3.4 |
| 128 | 7.08 | 8.82 | 0.092 | 12.31 | 21.28 | 0.6 |

**Table 14** The ablation study on filter number $F$ of the rotation module in MVMSR on CityStreet. Here $L = 3$ and $Q = 45$. G(1) and G(2) are GAME(1) and GAME(2), respectively.

| $L$ | MAE | MSE | NAE | G(1) | G(2) | FPS |
|---|---|---|---|---|---|---|
| 1 | 7.32 | 8.81 | 0.090 | 12.27 | 20.90 | **6.8** |
| 3 | 6.98 | 8.49 | 0.086 | **11.39** | **19.79** | 5.5 |
| 5 | **6.63** | **8.25** | **0.082** | 12.11 | 20.70 | 4.5 |
| 7 | 6.77 | 8.46 | 0.086 | 11.45 | 20.40 | 3.8 |
| 9 | 7.00 | 8.77 | 0.082 | 11.56 | 20.24 | 3.3 |

**Table 15** The ablation study on layer number $L$ of the rotation module in MVMSR on CityStreet. Here $F = 32$ and $Q = 45°$.

| $Q$ | MAE | MSE | NAE | G(1) | G(2) | FPS |
|---|---|---|---|---|---|---|
| 15° | 7.48 | 9.04 | 0.092 | 12.68 | 21.60 | 2.0 |
| 30° | 7.26 | 9.02 | 0.085 | **11.06** | **19.34** | 3.5 |
| 45° | **6.63** | **8.25** | **0.082** | 12.11 | 20.70 | 4.5 |
| 60° | 7.02 | 8.56 | 0.084 | 11.54 | 20.19 | 5.2 |
| 75° | 7.01 | 8.47 | 0.089 | 11.96 | 20.65 | 5.8 |
| 90° | 7.20 | 8.63 | 0.088 | 11.96 | 20.17 | **6.1** |

**Table 16** The ablation study on rotation quantization angle $Q$ of the rotation module in MVMSR on CityStreet. Here $F = 32$ and $L = 5$.

performance in terms of MAE, MSE and NAE metric while $L = 3$ achieves the best performance in terms of GAME(1) and GAME(2). We use $L = 5$ in the remaining experiments.

**Quantization angle $Q$:** Finally, we perform an ablation study on the quantization angle $Q$ in the rotation selection module of MVMSR on CityStreet, and the results are presented in Table 16. The choice of the quantization angle $Q$ involves the balance between the benefit of rotation selection and the extra learning complexity caused by the multi-rotations of the features. $Q = 45°$ has the best performance, compared to larger and smaller quantization angles on CityStreet.

In the main experiments tables (Table 3 and 9), we report the result on CityStreet with $F = 32$, $L = 3$ and $Q = 45°$ for better overall performance in terms of all evaluation metrics than the comparison methods.

### 6.3.5 Detection+ReID with Detection or ReID ground-truth

In crowd scenes, detection methods are limited by severe occlusions among the crowd, while ReID methods are hindered by detection errors, partial occlusions,
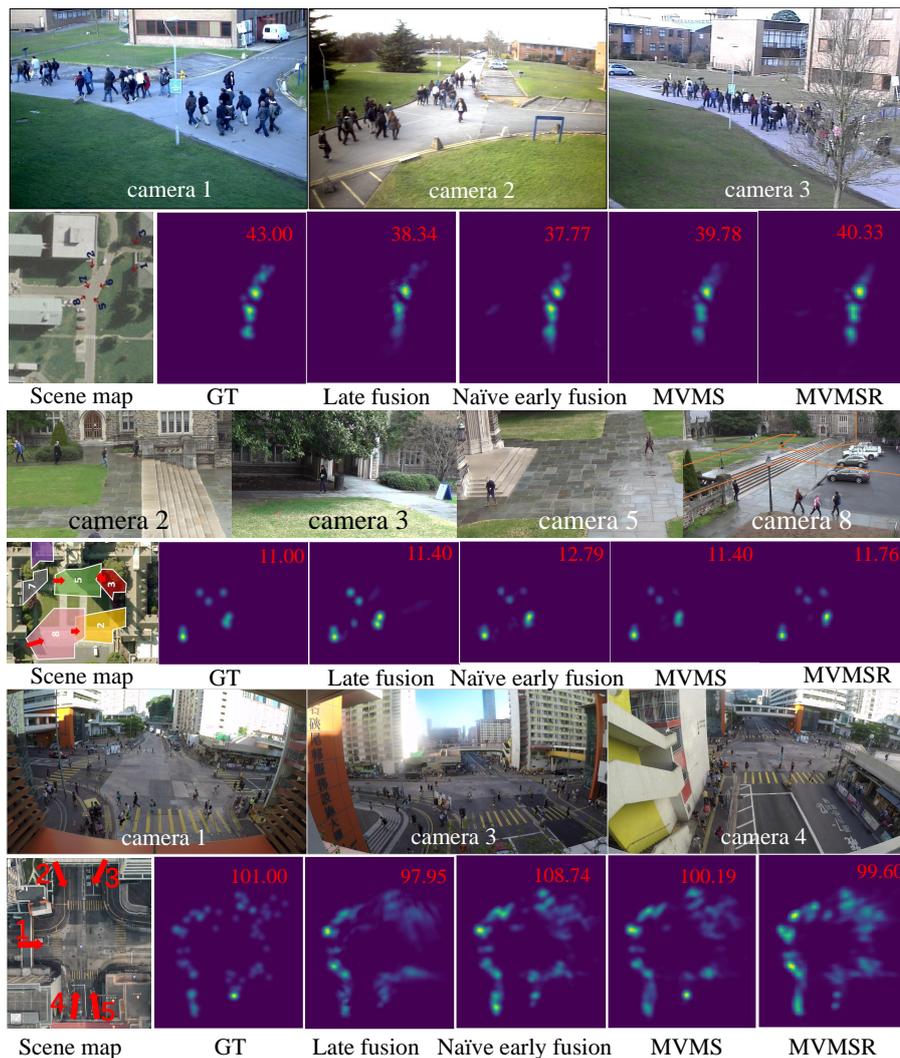
| Method | FPS |
|---|---|
| Dmap weighted (FCN-7) | 58.5 |
| Dmap weighted (CSR-net) | 9.4 |
| Detection+ReID | 0.3 |
| Feature concatenation | 3.4 |
| Stitching | 8.7 |
| Late fusion (FCN-7) | 27.8 |
| Naïve early fusion (FCN-7) | 30.7 |
| MVMS (FCN-7) | 19.8 |
| MVMSR (FCN-7) | 11.1 |
| Late fusion (CSR-net) | 5.2 |
| Naïve early fusion (CSR-net) | 8.5 |
| MVMS (CSR-net) | 7.9 |
| MVMSR (CSR-net) | 5.5 |
| Late fusion (LCC) | 3.9 |
| Naïve early fusion (LCC) | 5.0 |
| MVMS (LCC) | 4.4 |
| MVMSR (LCC) | 0.5 |

**Table 17** Running speed comparison on the CityStreet dataset.

scale changes between cameras, and low image-patch resolution. To illustrate the difficulties, we use the ground-truth inter-camera associations (*i.e.*, the best possible ReID) on the people detections and get counting MAE 30.3 on CityStreet, which is worse than our density-map fusion methods. Likewise, we apply ReID (Circle2021) on the ground-truth person boxes (*i.e.*, the best possible detector), and get counting MAE 10.7. Integrating multi-view detection and ReID for multi-view crowd counting would be interesting future work, and our dataset could serve as a test-bed.

### 6.3.6 Running speed comparison

We compare the running speed of different methods on the CityStreet dataset in the Table 17. Times are recorded for an Intel Xeon CPU E5-2543@3.30GHz with a Nvidia Geforce GTX 1080 Ti GPU. Generally, larger or deeper backbones (CSR-net and LCC) have lower running speed than lighter backbones (FCN-7). Comparing our methods, naïve early fusion is faster than late fusion because only one decoder module is needed to predict the scene-level density map, whereas late fusion has additional computations for predicting the density maps for each camera-view. MVMS is slower than naïve early fusion because MVMS extracts multi-scale features, compared to a single feature scale for naïve early fusion. Finally, MVMSR is slower than MVMS due to the additional network depth from the rotation module. The comparison method Dmap weighted (FCN-7) is faster than Dmap weighted (CSR-net) and Dmap_weighted (CSR-net), due its smaller backbone. For Detection+ReID (Circle2020), it uses large ReID models, so the running speeds is slow. The Stitching method's speed is similar to Dmap weighted (CSR-net),

**Fig. 14** The visualization results of the proposed multi-view fusion counting methods on PETS2009, DukeMTMC and CityStreet.

due to the same backbone model. For feature concatenation method, a large CNN model is used, so the running speed is comparable to our method's speed with larger backbones.

## 7 Discussion and Conclusion

In this paper, we propose a DNNs-based multi-view counting framework that fuses camera-views to predict scene-level ground-plane density maps for wide-area crowd counting. Both late fusion of density maps and early fusion of feature maps are studied. For late fusion, a projection normalization method is proposed to counter the effects of stretching caused by the projection operation. For early fusion, a multi-scale approach is proposed that selects features that have consistent

scales across views. We also propose a rotation selection module to handle rotated features introduced by the projection operation. To advance research in multi-view counting, we collect a new dataset of large scene containing a street intersection with large crowds. From the experiment results, our methods' performance gain over other comparison methods are larger on CityStreet, which is a larger and crowded scene. On the other hand, when the scene is not crowded enough, such as DukeMTMC, other methods can also achieve good performance. Nonetheless, our methods MVMS/MVMSR achieve the best performance on all 3 datasets.

In this paper, we focus on multi-camera counting when camera calibrations are known (like many other multi-camera vision tasks, such as 3D human pose estimation and multi-camera detection and tracking). The situation when the surveillance cameras orientation or

intrinsic parameters change gradually is interesting future work. One way to handle this situation would be to build an automatic calibration system, such as AutoClib (Bhardwaj et al. 2018), which uses car type and size to help calibrate the cameras, or (Ammar Abbas and Zisserman 2019), which computes a homography matrix for transforming the image to a geometrically correct bird's eye (overhead) view. Other calibration methods from 3D reconstruction (Agarwal et al. 2011; Snavely et al. 2006) could also be used.

Besides, adapting our framework to moving cameras and unknown camera parameters (using the full spatial transformer net) is interesting future work. In addition, we have trained and tested the network on each dataset individually. Another interesting future direction is on *cross-scene* multi-view counting, where the scenes in the test set are distinct from those in the training set – however, this requires more multi-view scenes to be collected. Since collecting videos of real scenes is difficult, especially during the pandemic, one recent work (Zhang et al. 2021) has collected a synthetic dataset for cross-scene cross-view multi-view counting.

# References

Agarwal S, Furukawa Y, Snavely N, Simon I, Curless B, Seitz SM, Szeliski R (2011) Building rome in a day. Communications of the ACM 54(10):105–112 21

Ammar Abbas S, Zisserman A (2019) A geometric approach to obtain a bird's eye view from an image. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops 21

Bhardwaj R, Tummala GK, Ramalingam G, Ramjee R, Sinha P (2018) Autocalib: Automatic traffic camera calibration at scale. ACM Transactions on Sensor Networks (TOSN) 14(3-4):1–27 21

Cao X, Wang Z, Zhao Y, Su F (2018) Scale aggregation network for accurate and efficient crowd counting. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 734–750 2, 4

Chan AB, Vasconcelos N (2012) Counting people with low-level features and bayesian regression. IEEE Transactions on Image Processing 21(4):2160–2177 3

Chan AB, Liang ZSJ, Vasconcelos N (2008) Privacy preserving crowd monitoring: Counting people without people models or tracking. In: Computer Vision and Pattern Recognition, pp 1–7 3

Chen C, Li G, Xu R, Chen T, Wang M, Lin L (2019) Clusternet: Deep hierarchical cluster network with rigorously rotation-invariant representation for point cloud analysis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 4994–5002 5

Chen CL, Chen K, Gong S, Xiang T (2013) Crowd Counting and Profiling: Methodology and Evaluation. Springer New York 3

Chen K, Chen LC, Gong S, Xiang T (2012) Feature mining for localised crowd counting. In: BMVC 3

Cheng Z, Qin L, Huang Q, Yan S, Tian Q (2014) Recognizing human group action by layered model with multiple cues. Neurocomputing 136:124–135 3

Cohen T, Welling M (2016) Group equivariant convolutional networks. In: International conference on machine learning, pp 2990–2999 5

Dieleman S, De Fauw J, Kavukcuoglu K (2016) Exploiting cyclic symmetry in convolutional neural networks. arXiv preprint arXiv:160202660 5

Dittrich F, de Oliveira LE, Britto Jr AS, Koerich AL (2017) People counting in crowded and outdoor scenes using a hybrid multi-camera approach. arXiv preprint arXiv:170400326 3, 4, 14, 15

Eiselein V, Fradi H, Keller I, Sikora T, Dugelay JL (2013) Enhancing human detection using crowd density measures and an adaptive correction filter. In: 10th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), IEEE, pp 19–24 2

Ferryman J, Shahrokni A (2009) Pets2009: Dataset and challenge. In: 2009 Twelfth IEEE International Workshop on Performance Evaluation of Tracking and Surveillance, IEEE, pp 1–6 3, 11, 14, 16

Gall J, Yao A, Razavi N, Van Gool L, Lempitsky V (2011) Hough forests for object detection, tracking, and action recognition. IEEE transactions on pattern analysis and machine intelligence 33(11):2188–2202 3

Gao H, Ji S (2017) Efficient and invariant convolutional neural networks for dense prediction. In: 2017 IEEE International Conference on Data Mining (ICDM), IEEE, pp 871–876 5

Ge W, Collins RT (2010) Crowd detection with a multiview sampler. In: European Conference on Computer Vision, pp 324–337 4

Guerrero-Gómez-Olmedo R, Torre-Jiménez B, López-Sastre R, Maldonado-Bascón S, Onoro-Rubio D (2015) Extremely overlapping vehicle counting. In: Iberian Conference on Pattern Recognition and Image Analysis, Springer, pp 423–431 13

Idrees H, Tayyab M, Athrey K, Zhang D, Al-Maadeed S, Rajpoot N, Shah M (2018) Composition loss for counting, density map estimation and localization in dense crowds. In: Proceedings of the European Conference on Computer Vision (ECCV) 2, 4

Jaderberg M, Simonyan K, Zisserman A, Kavukcuoglu K (2015) Spatial transformer networks. In: Advances in Neural Information Processing Systems (NIPS), pp 2017–2025 5, 6, 10

Jiang X, Xiao Z, Zhang B, Zhen X, Cao X, Doermann D, Shao L (2019) Crowd counting and density estimation by trellis encoder-decoder networks. In: CVPR, pp 6133–6142 4

Joachims T (1998) Text categorization with support vector machines: Learning with many relevant features. In: European conference on machine learning, Springer, pp 137–142 3

Junior JCSJ, Musse SR, Jung CR (2010) Crowd analysis using computer vision techniques. IEEE Signal Processing Magazine 27(5):66–77 3

Kang D, Chan A (2018) Crowd counting by adaptively fusing predictions from an image pyramid. In: BMVC 4

Kang D, Dhar D, Chan A (2017) Incorporating side information by adaptive convolution. In: Advances in Neural Information Processing Systems, pp 3867–3877 4

Kang D, Ma Z, Chan AB (2018) Beyond counting: comparisons of density maps for crowd analysis tasks-counting, detection, and tracking. IEEE Transactions on Circuits and Systems for Video Technology 2

Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems, pp 1097–1105 3

Laptev D, Savinov N, Buhmann JM, Pollefeys M (2016) Ti-pooling: transformation-invariant pooling for feature learning in convolutional neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 289–297 5

Leal-Taixé L, Milan A, Reid I, Roth S, Schindler K (2015) Motchallenge 2015: Towards a benchmark for multi-target tracking. arXiv preprint arXiv:150401942 11

Lempitsky V, Zisserman A (2010) Learning to count objects in images. In: Advances in Neural Information Processing Systems, pp 1324–1332 3

Li J, Huang L, Liu C (2012) People counting across multiple cameras for intelligent video surveillance. In: IEEE Ninth International Conference on Advanced Video and Signal-Based Surveillance (AVSS), IEEE, pp 178–183 3, 4

Li Y, Zhang X, Chen D (2018) Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 1091–1100 4, 6, 17

Lian D, Li J, Zheng J, Luo W, Gao S (2019) Density map regression guided detection network for rgb-d crowd counting and localization. In: CVPR, pp 1821–1830 4

Liu C, Weng X, Mu Y (2019a) Recurrent attentive zooming for joint crowd counting and precise localization. In: CVPR, pp 1217–1226 4

Liu J, Gao C, Meng D, Hauptmann AG (2018) Decidenet: Counting varying density crowds through attention guided detection and density estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 5197–5206 4

Liu L, Qiu Z, Li G, Liu S, Ouyang W, Lin L (2019b) Crowd counting with deep structured scale integration network. In: The IEEE International Conference on Computer Vision (ICCV) 4

Liu W, Salzmann M, Fua P (2019c) Context-aware crowd counting. In: CVPR, pp 5099–5108 4

Liu X, Yang J, Ding W, Wang T, Wang Z, Xiong J (2020) Adaptive mixture regression network with local counting map for crowd counting. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIV 16, Springer, pp 241–257 6, 13, 17

Ma H, Zeng C, Ling CX (2012) A reliable people counting system via multiple cameras. ACM Transactions on Intelligent Systems and Technology (TIST) 3(2):31 3, 4

Ma Z, Yu L, Chan AB (2015) Small instance detection by integer programming on object density maps. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3689–3697 2

Maddalena L, Petrosino A, Russo F (2014) People counting by learning their appearance in a multi-view camera environment. Pattern Recognition Letters 36:125–134 3, 4

Marana A, Costa LdF, Lotufo R, Velastin S (1998) On the efficacy of texture analysis for crowd monitoring. In: International Symposium on Computer Graphics, Image Processing, and Vision, IEEE, pp 354–361 3

Marcos D, Volpi M, Komodakis N, Tuia D (2017) Rotation equivariant vector field networks. In: Proceedings of the IEEE International Conference on Computer Vision, pp 5048–5057 5

Onoro-Rubio D, López-Sastre RJ (2016) Towards perspective-free object counting with deep learning. In: European Conference on Computer Vision, Springer, pp 615–629 4

Paragios N, Ramesh V (2001) A mrf-based approach for real-time subway monitoring. In: Computer Vision and Pattern Recognition, IEEE, vol 1 3

Pham VQ, Kozakaya T, Yamaguchi O, Okada R (2015) Count forest: Co-voting uncertain number of targets using random forest for crowd density estimation. In: Proceedings of the IEEE International Conference on Computer Vision, pp 3253–3261 3

Ranjan V, Le H, Hoai M (2018) Iterative crowd counting. In: ECCV, pp 270–285 4

Ren S, He K, Girshick R, Sun J (2015) Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in neural information processing systems, pp 91–99 13

Ren W, Kang D, Tang Y, Chan AB (2018) Fusing crowd density maps and visual object trackers for people tracking in crowd scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 5353–5362 2

Ristani E, Solera F, Zou R, Cucchiara R, Tomasi C (2016) Performance measures and a data set for multi-target, multi-camera tracking. In: European Conference on Computer Vision workshop on Benchmarking Multi-Target Tracking 3, 11, 14, 16

Rodriguez M, Laptev I, Sivic J, Audibert JY (2011) Density-aware person detection and tracking in crowds. In: IEEE International Conference on Computer Vision (ICCV), IEEE, pp 2423–2430 2

Ryan D, Denman S, Fookes C, Sridharan S (2014) Scene invariant multi camera crowd counting. Pattern Recognition Letters 44(8):98–112 4, 12, 13

Sabzmeydani P, Mori G (2007) Detecting pedestrians by learning shapelet features. In: IEEE Conference on Computer Vision and Pattern Recognition, IEEE, pp 1–8 3

Sam DB, Surya S, Babu RV (2017) Switching convolutional neural network for crowd counting. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, vol 1, p 6 4

Shen Z, Xu Y, Ni B, Wang M, Hu J, Yang X (2018) Crowd counting via adversarial cross-scale consistency pursuit. In: Computer Vision and Pattern Recognition, pp 5245–5254 4

Shi M, Yang Z, Xu C, Chen Q (2019) Revisiting perspective information for efficient crowd counting. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 7279–7288 4

Shi Z, Zhang L, Liu Y, Cao X, Ye Y, Cheng MM, Zheng G (2018) Crowd counting with deep negative correlation learning. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 5382–5390 4

Sindagi VA, Patel VM (2017) Generating high-quality crowd density maps using contextual pyramid cnns. In: IEEE International Conference on Computer Vision (ICCV), IEEE, pp 1879–1888 2, 4

Sindagi VA, Patel VM (2018) A survey of recent advances in cnn-based single image crowd counting and density esti-

mation. Pattern Recognition Letters 107:3–16 1, 3

Snavely N, Seitz SM, Szeliski R (2006) Photo tourism: exploring photo collections in 3d. In: ACM siggraph 2006 papers, pp 835–846 21

Sun Y, Cheng C, Zhang Y, Zhang C, Zheng L, Wang Z, Wei Y (2020) Circle loss: A unified perspective of pair similarity optimization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 6398–6407 13

Tang N, Lin YY, Weng MF, Liao HY (2014) Cross-camera knowledge transfer for multiview people counting. IEEE Transactions on Image Processing 24(1):80–93 4

Viola P, Jones MJ (2004) Robust real-time face detection. International journal of computer vision 57(2):137–154 3

Viola P, Jones MJ, Snow D (2005) Detecting pedestrians using patterns of motion and appearance. International Journal of Computer Vision 63(2):153–161 3

Wang G, Yuan Y, Chen X, Li J, Zhou X (2018) Learning discriminative features with multiple granularities for person re-identification. In: Proceedings of the 26th ACM international conference on Multimedia, pp 274–282 13

Wang Q, Gao J, et al (2019) Learning from synthetic data for crowd counting in the wild. In: CVPR, pp 8198–8207 4

Wang Y, Zou Y (2016) Fast visual object counting via example-based density estimation. In: IEEE International Conference on Image Processing (ICIP), IEEE, pp 3653–3657 3

Weiler M, Hamprecht FA, Storath M (2018) Learning steerable filters for rotation equivariant cnns. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 849–858 5

Worrall D, Brostow G (2018) Cubenet: Equivariance to 3d rotation and translation. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 567–584 5

Wu B, Nevatia R (2007) Detection and tracking of multiple, partially occluded humans by bayesian combination of edgelet based part detectors. International Journal of Computer Vision 75(2):247–266 3

Xu B, Qiu G (2016) Crowd density estimation based on rich features and random projection forest. In: IEEE Winter Conference on Applications of Computer Vision (WACV), IEEE, pp 1–8 3

Xu C, Qiu K, Fu J, Bai S, Xu Y, Bai X (2019) Learn to scale: Generating multipolar normalized density maps for crowd counting. In: The IEEE International Conference on Computer Vision (ICCV) 4

Yan Z, Yuan Y, Zuo W, Tan X, Wang Y, Wen S, Ding E (2019a) Perspective-guided convolution networks for crowd counting. In: The IEEE International Conference on Computer Vision (ICCV) 4

Yan Z, Yuan Y, Zuo W, Tan X, Wang Y, Wen S, Ding E (2019b) Perspective-guided convolution networks for crowd counting. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp 952–961 4

Yang Y, Li G, Wu Z, Su L, Huang Q, Sebe N (2020) Reverse perspective network for perspective-aware object counting. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 4374–4383 4

Zhang C, Li H, Wang X, Yang X (2015) Cross-scene crowd counting via deep convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 833–841 4

Zhang Q, Chan AB (2019) Wide-area crowd counting via ground-plane density maps and multi-view fusion cnns. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 8297–8306 4, 5

Zhang Q, Lin W, Chan AB (2021) Cross-view cross-scene multi-view crowd counting. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 557–567 15, 21

Zhang Y, Zhou D, Chen S, Gao S, Ma Y (2016) Single-image crowd counting via multi-column convolutional neural network. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 589–597 4

Zhang Z (2000) A flexible new technique for camera calibration. IEEE Transactions on pattern analysis and machine intelligence 22 11

Zheng L, Li Y, Mu Y (2021) Learning factorized cross-view fusion for multi-view crowd counting. In: 2021 IEEE International Conference on Multimedia and Expo (ICME), IEEE, pp 1–6 5, 15