# Single-Frame-Based Deep View Synchronization for Unsynchronized Multicamera Surveillance

Qi Zhang<sup>D</sup>, Student Member, IEEE, and Antoni B. Chan<sup>D</sup>, Senior Member, IEEE

Abstract-Multicamera surveillance has been an active 1 research topic for understanding and modeling scenes. Compared 2 to a single camera, multicameras provide larger field-of-view 3 and more object cues, and the related applications are multi-4 view counting, multiview tracking, 3-D pose estimation or 3-D 5 reconstruction, and so on. It is usually assumed that the cameras are all temporally synchronized when designing models for these 7 multicamera-based tasks. However, this assumption is not always 8 valid, especially for multicamera systems with network transmis-9 sion delay and low frame rates due to limited network bandwidth, 10 resulting in desynchronization of the captured frames across 11 cameras. To handle the issue of unsynchronized multicameras, 12 in this article, we propose a synchronization model that works 13 in conjunction with existing deep neural network (DNN)-based 14 multiview models, thus avoiding the redesign of the whole model. 15 We consider two variants of the model, based on where in the 16 pipeline the synchronization occurs, scene-level synchronization 17 and camera-level synchronization. The view synchronization 18 step and the task-specific view fusion and prediction step are 19 unified in the same framework and trained in an end-to-end 20 fashion. Our view synchronization models are applied to different 21 DNNs-based multicamera vision tasks under the unsynchronized 22 setting, including multiview counting and 3-D pose estimation, 23 and achieve good performance compared to baselines. 24

*Index Terms*—Crowd counting, deep learning, image match ing, pose estimation, surveillance.

27

## I. INTRODUCTION

▼OMPARED to single cameras, multicamera networks 28 → allow better understanding and modeling of the 3-D 29 world through more dense sampling of information in a 3-D 30 scene [1]. Multicamera based vision tasks have been a popular 31 research field, especially deep learning-related tasks, such as 3-32 D pose estimation from multiple 2-D observations [2], [3], 3-D 33 reconstruction [4], [5], multiview tracking [6]-[8], multiview 34 crowd counting [9], and re-identification (ReID) [10]-[14]. 35 Usually, it is assumed that the multicameras are temporally 36 synchronized when designing deep neural networks (DNNs) 37 models, i.e., all cameras capture images at the same time point. 38

Manuscript received July 22, 2021; revised January 4, 2022 and March 22, 2022; accepted April 15, 2022. This work was supported by Research Grant Council of Hong Kong Special Administrative Region (SAR), China, under Grant TR32-101/15-R and Grant CityU 11212518. (*Corresponding author: Qi Zhang.*)

Qi Zhang is with the Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ), Shenzhen, Guangdong 518060, China (e-mail: qzhang364-c@my.cityu.edu.hk).

Antoni B. Chan is with the Department of Computer Science, City University of Hong Kong, Hong Kong SAR, China (e-mail: abchan@cityu.edu.hk).

Color versions of one or more figures in this article are available at https://doi.org/10.1109/TNNLS.2022.3170642.

Digital Object Identifier 10.1109/TNNLS.2022.3170642

However, the synchronization assumption for multicamera 39 systems may not always be valid in *practical applications* due 40 to a variety of reasons, such as dropped camera frames due 41 to limited network bandwidth or system resources, network 42 transmission delays, and so on. Other examples of situations 43 where camera synchronization is not possible include: 1) 44 using images captured from different camera systems; 2) using 45 images from social media to reconstruct the crowd at an event; 46 and 3) performing 3-D reconstruction of a dynamic scene 47 using video from a drone. Thus, handling unsynchronized mul-48 ticameras is an important issue in the adoption and practical 49 usage of multiview computer vision. 50

There are several possible methods to fix the prob-51 lem of unsynchronized cameras. The first method is using 52 hardware-based solutions to synchronize the capture times, 53 such as improving network bandwidth, or by using a central 54 clock to synchronize capture of all cameras in the multicamera 55 network. However, this will increase the cost and overhead 56 of the system, and is not possible when there is limited 57 bandwidth. The second method is to capture image sequences 58 from each camera, and then synchronize the images afterward 59 by determining the frame offset between cameras. The fineness 60 of the synchronization depends on the frame rate of the 61 image sequences. However, this method is not effective when 62 acquiring high frame-rate image sequences is not possible 63 due to limited the bandwidth and storage space, or the frame 64 latency between multicameras is random. The final method 65 is to modify the multiview model to handle unsynchronized 66 images, especially for low-frame-rate multicamera systems or 67 random frame latency between multicameras, such as introduc-68 ing new assumptions or relaxing the original constraints under 69 the unsynchronized setting. Existing approaches for handling 70 unsynchronized multicameras are largely based on optimiza-71 tion frameworks [15], [16], but are not directly applicable 72 to DNNs-based multiview methods, which have seen recent 73 successes in tracking [6], [7], 3-D pose estimation [2], and 74 crowd counting [9], [17]. 75

In this article, we propose a synchronization model that 76 operates in conjunction with existing DNN-based multiview 77 models by using single frames from each camera to deal 78 with low-frame-rate unsynchronized multicamera systems or 79 random frame latency between multicameras. Our proposed 80 model first synchronizes other views to a reference view using 81 a differentiable module, and then the synchronized multiviews 82 features are fused and decoded to obtain the task-oriented 83 output. As illustrated in Fig. 1, the synchronization can either 84

2162-237X © 2022 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.



Fig. 1. Two variants of the main pipeline for unsynchronized multiview prediction tasks: (top) SLS is performed after the projection on the scene-level feature representations and (bottom) CLS is performed on the camera-view feature maps before projection.

occur after the camera-to-scene (2-D-to-3-D) projection [Fig. 1 85 (top)] or before the projection [Fig. 1 (bottom)]. Thus, to fully 86 explore these options, we consider two variants of our model 87 that perform synchronization at different stages in the pipeline 88 (see Fig. 2): 1) scene-level synchronization (SLS) performs 89 the synchronization after projecting the camera features to 90 their 3-D scene representation and 2) camera-level synchro-91 nization (CLS) performs the synchronization between camera 92 views first, and then projects the synchronized 2-D feature 93 maps to their 3-D representations. In both cases, motion flow 94 between the cameras' feature maps are estimated and then 95 used to warp the feature maps to align with the reference 96 view (either at the scene-level or the camera-level). With 97 both variants, the view synchronization and the multiview 98 fusion are unified in the same framework and trained in 99 an end-to-end fashion. In this way, the original DNN-based 100 multiview model can be adapted to work in the unsynchronized 101 setting by adding the view synchronization module, thus 102 avoiding the need to design a new model. Furthermore, the 103 synchronization module only relies on content-based image 104 matching and camera geometry, and thus is widely applicable 105 to many DNNs-based multiview tasks, such as crowd counting, 106 tracking, 3-D pose estimation, and 3-D reconstruction. 107

<sup>108</sup> In summary, the contributions of this article are threefold.

- We propose an end-to-end trainable framework to handle the issue of unsynchronized multicamera images in DNNs-based multicamera vision tasks. To the best of our knowledge, this is the first study on DNNs-based *single-frame* synchronization of multiview cameras.
- 2) We propose two synchronization modules, SLS and camera-view level synchronization, which are based on image-based content matching that is guided by epipolar geometry. The synchronization modules can be applied to many different DNNs-based multiview tasks.
- 3) We conduct experiments on multiview counting and 3-D
   pose estimation from unsynchronized images, demon strating the efficacy of our approach.

The remainder of this article is organized as follows. We review related works in Section II. In Section III, we propose our single-frame multicamera synchronization methods, and in Section IV we present experiments on two applications, multiview crowd counting, and multiview 3-D human pose estimation. Finally, Section V concludes the 127 article.

In this section, we review DNN-based methods on synchronized multiview images and unsynchronized multiview video tasks, as well as traditional multiview video synchronization methods. We then review DNN-based image matching and flow estimation methods.

## A. DNN-Based Synchronized Multicamera Tasks

II.

Multicamera surveillance based on DNNs has been an 136 active research area. By utilizing multiview cues and the 137 strong mapping power of DNNs, many DNNs models have 138 been proposed to solve multiview surveillance tasks, such 139 as multiview tracking and detection [6], [7], [18], crowd 140 counting [9], 3-D reconstruction [4], [5], [19], [20], and 3-D 141 human pose estimation [2], [21]-[24]. Kar et al. [4] proposed 142 a deep learning 3-D reconstruction framework with differen-143 tiable feature projection and unprojection steps. Ye et al. [10] 144 proposed the collaboration ensemble learning for ReID with 145 middle-level sharable two-stream network. Iskakov et al. [2] 146 proposed volumetric aggregation of feature maps for 3-D 147 pose estimation. The DNN pipelines used for these mul-148 ticamera tasks can be generally divided into three stages: 149 the single-view feature extraction stage, the multiview fusion 150 stage to obtain a scene-level representation, and prediction 151 stage. Furthermore, all these DNN-based methods assume 152 that the input multiviews are synchronized, which is not 153 always possible in real multicamera surveillance systems, or in 154 multiview data from disparate sources (e.g., crowd-sourced 155 images). Therefore, relaxing the synchronization assumption 156 can allow more practical applications of multicamera vision 157 tasks in real world. 158

## B. Tasks on Unsynchronized Multicamera Video

Only a few works have considered computer vision tasks on 160 unsynchronized multicamera video. Zheng et al. [15] posed 161 the estimation of 3-D structure observed by multiple unsyn-162 chronized video cameras as the problem of dictionary learning. 163 Zhang et al. [16] proposed a multicamera motion segmenta-164 tion method using unsynchronized videos by combining shape 165 and dynamical information. Takahashi et al. [25] proposed a 166 method of estimating 3-D human pose from multiview videos 167 captured by unsynchronized and uncalibrated cameras by 168 utilizing the projections of joint as the corresponding points. 169 Albl et al. [26] presented a method for simultaneously esti-170 mating camera geometry and time shift from video sequences 171 from multiple unsynchronized cameras using minimal cor-172 respondence sets. Kuo et al. [27] addressed the problem of 173 aligning unsynchronized camera views by low and/or variable 174 frame rates using the intersections of corresponding object 175 trajectories to match views. 176

Note that all these methods assume that videos or image sequences are available to perform the synchronization. In contrast, our framework, which is motivated by practical low-fps systems, is solving a harder problem, where *only a single* 180

159



Fig. 2. General multiview pipeline consists of several stages: *camera-view feature extraction, feature projection, multiview feature fusion to obtain a scene-level representation, and prediction.* (a) **SLS** performs the synchronization *after the projection.* The unsynchronized projected features from the reference view and other views are *concatenated to predict the motion flow*, which is then used to *warp the other views' projected features to match those of the reference view*. (b) **CLS** performs the synchronization *before the projection.* The unsynchronized camera-view features from the reference view and other views are *matched together to predict the motion flow*, which is used to *warp features from other camera views to the reference view.* 

*image* is available from each camera view, i.e., there is no 181 temporal information available. Furthermore, these methods 182 pose frame synchronization as optimization problems that are 183 applicable only to the particular multiview task, and cannot be 184 directly applied to DNN-based multiview models. In contrast, 185 we propose a synchronization module that can be broadly 186 applied to many DNN-based multicamera models, enabling 187 their use with unsynchronized inputs. 188

## 189 C. Traditional Methods for Multiview Video Synchronization

Traditional synchronization methods usually serve as 190 a preprocessing step for multicamera surveillance tasks. 191 Except audio-based synchronization like [28], most tradi-192 tional camera synchronization methods rely on videos or 193 image sequences and hand-crafted features for camera align-194 ment/synchronization [29]-[33]. Typical approaches recover 195 the temporal offset by matching features extracted from the 196 videos, e.g., space-time feature trajectories [34]-[36], image 197 features [37], low-level temporal signals based on funda-198 mental matrices [38], silhouette motion [39], and relative 199 object motion [40]. The accuracy of feature matching is 200 improved using epipolar geometry [37], [39] and rank con-201 straints [35]. Caspi et al. [34] proposed to use the space-time 202 feature trajectories matching instead of feature-points match-203 ing to reduce the search space. Dai et al. [29] proposed an 204 iterative procedure to achieve the alignment in space and 205 time with the homography assumption in spatial domain. 206 Imre and Hilton [37] utilized image feature correspondences 207 and epipolar geometry to find the corresponding frame indices 208 and computes the relative frame rate and offset by fitting 209 a 2-D line to the index correspondences. Meyer et al. [36] 210

estimated the frame accurate offset by analyzing the tra-21 jectories and matching their characteristic time patterns. 212 Pundik and Moses [38] presented a method for online syn-213 chronization that relied on the video sequences with known 214 fundamental matrix to compute low-level temporal signals 215 for matching. Rao et al. [35] proposed the rank constraint of 216 corresponding points in two views to measure the similarity 217 between trajectories to avoid the noise sensitivity of the 218 fundamental matrix. Sinha and Pollefeys [39] proposed a 219 Random sample consensus (RANSAC)-based algorithm that 220 computed the epipolar geometry and synchronization of a pair 221 of cameras from the motion of silhouettes in videos. Tresadern 222 and Reid [32] estimated possible synchronization parameters 223 via the Hough transform and refined these parameters using 224 nonlinear optimization methods. Yan and Pollefeys [33] relied 225 on correlating space-time interest point distribution in time 226 between videos which represented events in video that had 227 high variation in both space and time. Gaspar et al. [40] syn-228 chronized two independently moving cameras via the relative 229 motion between objects and known camera intrinsic. 230

The main disadvantages for these traditional camera synchronization methods are.

- 1) Videos and image sequences are required, which might not be available in practical multicamera systems with limited network bandwidth and storage.
- 2) A fixed frame rate of the multicameras are usually assumed, which means random frame dropping cannot be handled (except [38]).
- Feature matching is based on hand-crafted features, which lack representation ability, or known image correspondences, which requires extra manual annotations and may not always be available.

231

232

233

234

235

236

237

Compared with these methods, we consider a more practical 243 and difficult setting: only single-frames and no videos (no tem-244 poral information) are available, which means that these tra-245 ditional video-based methods are not suitable solutions. These 246 traditional methods perform image content matching using 247 hand-crafted features and traditional matching algorithms, 248 while in contrast, our method uses DNN-based image match-249 ing. Because we also assume that only single-frames are avail-250 able, our method also requires DNN-based motion estimation 251 to estimate a frame's features after synchronization. Finally, 252 our synchronization module is end-to-end trainable with exist-253 ing multiview DNNs and thus avoids the redesign of the whole 254 DNNs models to handle unsynchronized multicameras. 255

## <sup>256</sup> D. DNN-Based Image Matching and Flow Estimation

Image matching and optical flow estimation both involve 257 estimating image-to-image correspondences, which is related 258 to frame synchronization of multiviews. We mainly review 259 the DNN-based image matching [41]-[43] or optical flow 260 estimation methods [44]-[46], which inspire us to solve the 261 unsynchronized multicamera based problems in a DNN-based 262 fashion. DNN flow [47] proposed an image matching method 263 based on a DNN feature pyramid in a coarse-to-fine optimiza-264 tion manner. FlowNet [48] predicted the optical flow from 265 DNNs with feature concatenation and correlation. SpyNet [49] 266 combined a classical spatial-pyramid formulation with deep 267 learning and estimated large motions in a coarse-to-fine 268 approach by warping one image to the other at each pyramid 269 level by the current flow estimate and computing an update 270 to the flow. Rocco et al. [41] addressed image correspondence 271 problem using a convolutional neural network architecture that 272 mimics classic image matching algorithms. PWC-Net [50] 273 uses a feature pyramid and one image feature map is warped 274 to the other at each scale, which is guided by the upsampled 275 optical flow estimated from the previous scale. Lai et al. [51] 276 proposed a single network to jointly learn spatiotemporal 277 correspondence for stereo matching and flow estimation. 278

Our method is related to the DNN-based image matching and optical flow estimation, but the difference is still significant.

- Typical image/geometric matching only involves *either* a camera view angle transformation (e.g., camera relative pose estimation, stereo matching) *or* a small time change in the same view (optical flow estimation), while *both* factors appear in our problem, which makes our problem harder.
- 288
   2) Image/geometric matching is directly supervised by the correspondence of two images, while the multiview fusion ground-truth in the 3-D world is used as supervisory signal in our problem.
- 3) The 2-D-to-3-D projection causes ambiguity for multiview feature fusion, which also causes difficulties for view synchronization.

## 295 III. SINGLE-FRAME DNNS MULTICAMERA 296 SYNCHRONIZATION

In this section, we propose our single-frame synchronization model for DNN-based multiview models. The temporal offset

between cameras can either be constant latency for each 299 camera (the same offset over time), or random latency (random 300 offsets over time). Similar to most multiview methods [2], [7], 301 [17], [20], we assume that the cameras are static and the cam-302 eras' intrinsic and extrinsic parameters are known. The main 303 idea of our method is to choose a camera view as the 304 reference view, and then use the view synchronization model 305 to warp the other camera views to be synchronized with the 306 reference view. The synchronization model should be general 307 enough to handle both constant and random latencies between 308 cameras, in order to work under various conditions causing 309 desynchronization. 310

DNNs models for the multicamera surveillance tasks typ-311 ically consist of three stages (see in Fig. 1): Single-view 312 feature extraction, which extracts single-view features of 313 the input camera views. Multiview feature projection and 314 fusion, where a fixed differentiable projection layer is first 315 adopted to project the single-view features to the 3-D coordi-316 nate map and then the projected multiview features are fused 317 together to form the scene-level representation. The projec-318 tion layer depends the application task, and our framework 319 can generally handle any differentiable projection layer. For 320 example, for multiview counting [9], the projection maps 321 the 2-D camera view to the 3-D scene plane at the average 322 person height (assuming all camera pixels fall on the same 323 height plane), while for 3-D pose estimation [2], the projection 324 copies features along a view-ray in the 3-D grid, assuming 325 an unknown height of each camera-view pixel. Prediction, 326 where the decoder predicts the final result in the 3-D coor-327 dinate map, such as ground-plane density maps [9] or 3-D 328 reconstruction [4]. 329

In Fig. 2, we take multiview crowd counting [9] as an 330 example to show the pipeline of the proposed single-frame-331 based view synchronization model. In the multiview fusion 332 model, we denote the input multiview frames as  $\{I_i^{t_0}\}_{i=0}^{n-1}$ , 333 where i denotes the camera view id and n is the input camera 334 view number, and superscript  $t_0$  indicates that the frames 335 are all captured at the same time point  $t_0$ , corresponding 336 to the synchronized multicamera setup. After being fed into 337 the single-view feature extractor F, the extracted features are 338 denoted as 339

$$F_i^{t_0} = F(I_i^{t_0}), \quad i \in \{0, 1, \dots, n-1\}.$$
 (1) 340

For multiview counting [9], the projection  $\mathcal{P}$  maps the 2-D <sup>341</sup> camera view to the 3-D scene plane at the average person <sup>342</sup> height. After projection layer  $\mathcal{P}$ , the projected multiview <sup>343</sup> features are <sup>344</sup>

$$\mathcal{F}_{i}^{t_{0}} = \mathcal{P}(F_{i}^{t_{0}}), \quad i \in \{0, 1, \dots, n-1\}.$$
 (2) 345

We use U to denote the fusion operation (e.g., concatenation and max-pooling) of the projected multiview features, thus the fused feature is  $U(\mathcal{F}_0^{t_0}, \ldots, \mathcal{F}_{n-1}^{t_0})$ . Finally, the decoder D is applied to obtain the final prediction  $V_p$  349

$$V_p = D(U(\mathcal{F}_0^{t_0}, \dots, \mathcal{F}_{n-1}^{t_0}))$$
 350

$$= D(U(\mathcal{P}(F_0^{t_0}), \dots, \mathcal{P}(F_{n-1}^{t_0}))).$$
(3) 35

However, when the input multicameras frames are not synchronized, denoted as  $\{I_i^{t_i}\}_{i=0}^{n-1}$ , the capture time for the *i*th 353

The view synchronization model can be embedded into one 356 of the first two stages, synchronizing the extracted single-view 357 features  $\{F_i^{t_i}\}$  or projected features  $\{\mathcal{F}_i^{t_i}\}$ , without the need to 358 redesign a new architecture. Thus, we propose two variants 359 of the synchronization model: 1) SLS, where the projected 360 features  $\{\mathcal{F}_i^{t_i}\}$  from different camera views are synchronized 361 during multicamera feature fusion and 2) CLS, where the 362 camera view features  $\{F_i^{t_i}\}$  are synchronized before projection 363 and fusion. We present the details of the two synchronization 364 models next. Note that we first consider the case when 365 both synchronized and unsynchronized multiview images are 366 available for training (but not available in the testing stage). 367 We then extend this to the case when only unsynchronized 368 training images are available. 369

## 370 A. Scene-Level Synchronization

SLS works by synchronizing the multicamera features after the projection stage in the multiview pipeline. The pipeline for SLS is shown in Fig. 2(a).

1) Synchronization Module: Without loss in generality, we choose one view (denoted as view 0) as the reference view, and other views are to be synchronized to this reference view. We first assume that synchronized frame pairs are available in the training stage. The frames are  $I_0^{t_0}$  from reference view 0 captured at reference time  $t_0$ , and  $I_i^{t_0}$  and  $I_i^{t_i}$  from view *i* ( $i \in \{1, 2, ..., n-1\}$ ) taken at times  $t_0$  and  $t_i$ . Note that frames  $(I_0^{t_0}, I_i^{t_0})$  are synchronized, while  $(I_0^{t_0}, I_i^{t_i})$  are not.

The synchronization module consists of the following 382 stages. First, camera frame feature maps  $(F_0^{t_0}, F_i^{t_0}, F_i^{t_i})$  (both 383 synchronized and unsynchronized frames) are extracted and 384 projected to the 3-D world space, resulting in the projected 385 feature maps  $(\mathcal{F}_0^{t_0}, \mathcal{F}_i^{t_0}, \mathcal{F}_i^{t_i})$ . Second, synchronization is per-386 formed between the reference view 0 and each other view 387 *i*. The projected feature map  $\mathcal{F}_0^{t_0}$  from the reference view is 388 concatenated with the projected feature map  $\mathcal{F}_{i}^{t_{i}}$  from view *i*, 389 and then fed into a motion flow estimation network  $\mathcal{M}_s$  to 390 predict the scene-level motion flow  $w_i$  between view *i* at time 391  $t_i$  and the reference view at time  $t_0$ 392

393 
$$w_i = \mathcal{M}_s(\operatorname{Cat}(\mathcal{F}_0^{t_0}, \mathcal{F}_i^{t_i})), \quad i \in \{1, \dots, n-1\}$$
 (4)

where Cat is the concatenation operation. The  $\mathcal{F}_{i}^{t_{i}}$  from view *i* is then synchronized with the reference view at time  $t_{0}$  using a warping transformation  $\mathcal{W}$  guided by  $w_{i}, \mathcal{W}(w_{i}, \mathcal{F}_{i}^{t_{i}})$ 

397

$$\hat{\mathcal{F}}_i^{t_0} = \mathcal{W}\big(w_i, \mathcal{F}_i^{t_i}\big), \quad i \in \{1, \dots, n-1\}$$
(5)

where  $\hat{\mathcal{F}}_{i}^{t_{0}}$  are the warped projected features of the *i*th view synchronized to time  $t_{0}$ . Note that the warping  $\mathcal{W}$  only applies spatial shifting to the feature map  $\mathcal{F}_{i}^{t_{i}}$ , i.e., it only changes the feature locations and does not change the feature values. Finally, the reference view features  $\mathcal{F}_{0}^{t_{0}}$  and estimated warped features of the other views  $\{\hat{\mathcal{F}}_{i}^{t_{0}}\}$  are fused and decoded to obtain the final scene-level prediction  $V_{p}$ 

405 
$$V_p = D(U(\mathcal{F}_0^{t_0}, \hat{\mathcal{F}}_1^{t_0}, \dots, \hat{\mathcal{F}}_{n-1}^{t_0}))$$
 (6)

$$= D(U(\mathcal{F}_{0}^{i_{0}}, \mathcal{W}(w_{1}, \mathcal{F}_{1}^{i_{1}}), \dots, \mathcal{W}(w_{n-1}, \mathcal{F}_{n-1}^{i_{n-1}}))).$$
(7)

In the testing stage, only unsynchronized frames  $(I_0^{t_0}, I_i^{t_i})$  are 407 available and the forward operations related to frame  $I_i^{t_0}$  are 408 removed from the network. 409

2) Training Loss: Two losses are used in the training stage. 410 The first loss is a task-specific prediction loss  $\ell_p$  between 411 the scene-level prediction  $V_p$  and the ground-truth  $V_{gt}$ . For 412 example, for multiview crowd counting  $\ell_p$  is the mean-square 413 error, and  $V_p$ ,  $V_{gt}$  are the predicted and ground-truth scene-414 level density maps. The second loss is on the multiview 415 feature synchronization in the multiview fusion stage. Since 416 the synced frame pairs are available during training, the feature 417 warping loss  $\ell_{\mathcal{W}}$  encourages the warped features to be similar 418 to the features of the original synced frame of view *i* 419

$$\ell_{\mathcal{W}}(w_i, \mathcal{F}_i^{t_0}, \mathcal{F}_i^{t_i}) = \operatorname{mse}(\mathcal{F}_i^{t_0}, \hat{\mathcal{F}}_i^{t_0})$$

$$\operatorname{mse}(\mathcal{T}_i^{t_0}, \mathcal{W}_i^{t_i})$$

$$(8)$$

$$= \operatorname{mse}(\mathcal{F}_{i}^{0}, \mathcal{W}(w_{i}, \mathcal{F}_{i}^{0}))$$
(8) 42

where mse is the mean-square error loss. Note that the warping 422  $\mathcal{W}$  only applies spatial shifting, and thus the minimization 423 of the warping loss  $\ell_{\mathcal{W}}$  in (8) will be based on the feature 424 alignment via scene-level motion flow  $w_i$  and not global 425 feature value changes (e.g., color correction). Finally, the 426 training loss combines the task loss and the warping loss 427 summed over all nonreference views 428

$$\ell = \ell_p \left( V_p, V_{gt} \right) + \gamma \sum_{i=1}^{n-1} \ell_{\mathcal{W}} \left( w_i, \mathcal{F}_i^{t_i}, \mathcal{F}_i^{t_i} \right)$$
(9) 429

where  $\gamma$  is a hyperparameter.

## B. Camera View-Level Synchronization

Each image pixels' height in 3-D space is unknown, and 432 thus the projection operation of multicamera DNNs mod-433 els [2], [9], [17] will either project each pixel to the same 434 assumed height level [9] (causing distortion when the true 435 pixel height is different), or to multiple height levels [2], [17] 436 (duplicating features along the view ray). These projection 437 cause the features to stretch along the view ray in the 3-D 438 scene, which makes their synchronization more difficult due to 439 their imprecise (stretched) and ambiguous (duplicated) nature. 440 Therefore, we also consider synchronization between camera 441 view features before the projection. The pipeline for CLS is 442 presented in Fig. 2(b). 443

1) Synchronization Model: The view synchronization model 444 is applied to each view separately. The camera view features 445  $(F_0^{t_0}, F_i^{t_i})$  from the unsynchronized reference view and view 446 *i* are first passed through a matching module (see below) 447 and then fed into the motion flow estimation network  $\mathcal{M}_{c}$ 448 to predict the camera-view motion flow  $w_i$  for view i. The 449 warping transformation W guided by  $w_i$  then warps the 450 camera-view features  $F_i^{t_i}$  from view *i* to be synchronized with 451 the reference view at time  $t_0$ 452

$$\hat{F}_i^{t_0} = W(w_i, F_i^{t_i}), \quad i \in \{1, \dots, n-1\}$$
 (10) 453

where  $\hat{F}_{i}^{t_{0}}$  is the warped camera-view features of view *i* 454 captured at time  $t_{i}$ , which is synchronized to reference view 455 0 captured at time  $t_{0}$ . Finally, the reference and warped camera 456 views are projected 457

$$\mathcal{F}_{0}^{t_{0}} = \mathcal{P}(F_{0}^{t_{0}}), \hat{\mathcal{F}}_{i}^{t_{0}} = \mathcal{P}(\hat{F}_{i}^{t_{0}}), \quad i \in \{1, \dots, n-1\}$$
 (11) 458



Fig. 3. Epipolar-guided weights. (a) In the synchronized setting, given the point (x, y) in view 0, the matched point (x', y') in view *i* must be on the epipolar line  $l_{xy}$ . (b) In the unsynchronized setting, we assume a Gaussian motion model of the matched feature location from time  $t_0$  to  $t_i$ . (c) Epipolar-guided weight mask is use to bias the feature matching toward high-probability regions according to the motion model.

and then fused and decoded to obtain the scene-level prediction  $V_p$ 

$$V_p = D(U(\mathcal{F}_0^{t_0}, \hat{\mathcal{F}}_1^{t_0}, \dots, \hat{\mathcal{F}}_{n-1}^{t_0}))$$
(12)

$$= D(U(\mathcal{P}(F_0^{t_0}), \mathcal{P}(\hat{F}_1^{t_0}), \dots, \mathcal{P}(\hat{F}_{n-1}^{t_0}))).$$
(13)

In the testing stage, only unsynchronized frames  $(I_0^{t_0}, I_i^{t_i})$  are available and the forward operations related to frame  $I_i^{t_0}$  are removed from the network.

<sup>466</sup> 2) *Matching Module:* We propose three methods to match features to predict the view-level motion flow. The first method concatenates the features  $(F_0^{t_0}, F_i^{t_i})$  and then feeds them into the motion flow estimation network  $\mathcal{M}_c$  to predict the motion flow  $w_i$ 

$$w_i = \mathcal{M}_c(\operatorname{Cat}(F_0^{t_0}, F_i^{t_i})), \quad i \in \{1, \dots, n-1\}.$$
(14)

The second method builds a correlation map  $C_i$  between features from each pair of spatial locations in  $F_0^{t_0}$  and  $F_i^{t_i}$ 

474 
$$C_i((x, y), (x', y')) = F_0^{t_0}(x, y)^T F_i^{t_i}(x', y')$$
(15)

which is then fed into the motion flow estimation network  $\mathcal{M}_c$ to predict the motion flow  $w_i$ 

477

471

 $w_i = \mathcal{M}_c(C_i), \quad i \in \{1, \dots, n-1\}.$  (16)

The third method incorporates camera geometry information 478 into the correlation map to suppress false matches. If both 479 cameras are synchronized at  $t_0$ , then according the multiview 480 geometry, each spatial location in view 0 must match a location 481 482 in view *i* on its corresponding epipolar line [Fig. 3(a)]. Thus, in the synchronized setting, detected matches that are not on 483 the epipolar line can be rejected as false matches. For our 484 unsynchronized setting, the matched location in view *i* remains 485 on the epipolar line only when its corresponding feature/object 486 does not move between times  $t_0$  and  $t_i$ . To handle the case 487 where the feature moves, we assume that a matched feature 488 in view *i* moves according to a Gaussian motion model with 489 standard deviation  $\sigma$  [Fig. 3(b)]. With the epipolar line and 490 motion model, we then build a weighting mask, with high 491 weights on locations with high probability of containing the 492 matched features, and vice versa. Specifically, we set the 493 mask  $M_i((x, y), (x', y')) = 1$  if (x', y') is on the epipolar 494 line induced by (x, y), and 0 otherwise, and then convolve 495 it with a 2-D Gaussian with standard deviation  $\sigma$  [Fig. 3(c)]. 496 We then apply the weight mask  $M_i$  on the correlation map 497  $\tilde{C}_i = M_i \odot C_i$ , which will suppress false matches that are not 498



Fig. 4. Multiscale estimation of motion flow.

consistent with the scene and motion model. Thus, the motion 499 flow  $w_i$  is 500

$$w_i = \mathcal{M}_c(\tilde{C}_i) = \mathcal{M}_c(M_i \odot C_i), \quad i \in \{1, \dots, n-1\}.$$
(17) 50

3) Multiscale Architecture: Multiscale feature extractors are 502 used in multicamera tasks like crowd counting [9] or to refine 503 the final prediction via multiscale prediction fusion [50], [51]. 504 Therefore, we next show how to incorporate multiscale feature 505 extractors with our CLS model.<sup>1</sup> Instead of performing the 506 view synchronization in each scale separately, the motion flow 507 estimate of neighbor scales is fused to refine the current scale's 508 estimate (see Fig. 4). In particular, let there be *m* scales in the 509 multiscale architecture and j denotes one scale in the scale 510 range  $\{1, 2, \ldots, m\}$ , with m the largest scale. The multiscale 511 predicted motion flow are fused as follows. 512

- 1) When j = 1 (the smallest scale), the correlation map  $C_i^{(1)}$  of scale 1 is fed into the motion flow estimation net to predict the motion flow  $w_i^{(1)}$  for scale 1.
- 2) For scales j > 1, first the difference between the correlation map  $C_i^{(j)}$  and the upsampled correlation map of the previous scale up $(C_i^{(j-1)})$  is fed into the motion flow estimation net to predict the residual of the motion flow between two scales, denoted as  $\tilde{w}_i^{(j)}$ .

3) The refined motion flow of scale j is

$$w_i^{(j)} = up\left(w_i^{(j-1)}\right) + \tilde{w}_i^{(j)}.$$
 (18) 522

521

4) *Training Loss:* Similar to SLS, a combination of two losses (scene-level prediction and feature synchronization) is used in the training stage. The scene-level prediction loss is the same as before. The feature synchronization loss encourages the warped camera-view features at each scale to match the features of the original synchronized frame

$$\ell_W = \mathrm{mse}\Big(F_i^{t_0,(j)}, \hat{F}_i^{t_0,(j)}\Big)$$
(19) 524

$$= \operatorname{mse}\left(F_{i}^{t_{0},(j)}, W\left(w_{i}^{(j)}, F_{i}^{t_{i},(j)}\right)\right).$$
(20) 530

Similar to SLS, the warping function W only applies spatial shifting, and thus the minimization of  $\ell_W$  in (20) will be based on feature alignment rather than feature value changes. Finally, the training loss is the combination of the prediction loss and the synchronization loss summed over all nonreference views 535

<sup>1</sup>No extra steps are needed to incorporate multiscale features with SLS because the synchronization occurs after the feature projection.

536 and scales

548

557

537 
$$\ell = \ell_p (V_p, V_{gt}) + \gamma \sum_{i=1}^{n-1} \sum_{j=1}^m \ell_W \left( w_i^{(j)}, \mathcal{F}_i^{t_0,(j)}, \mathcal{F}_i^{t_i,(j)} \right) \quad (21)$$

538 where  $\gamma$  is a hyperparameter.

## 539 C. Training With Only Unsynchronized Frames

In the previous models, we assume that both synchronized 540 and unsynchronized multicamera frames are available during 541 training. For more practical applications, we also consider 542 the case when only unsynchronized multiview frames are 543 available for training. In this case, for the SLS, the warping 544 feature loss  $\ell_W$  is replaced with a similarity loss  $\ell_s$  on the 545 projected features, to indirectly encourage synchronization of 546 the projected multiview features 547

$$\ell_s = \operatorname{mean}(1 - \cos(\mathcal{F}_0^{t_0}, \mathcal{W}(w_i, \mathcal{F}_i^{t_i})))$$
(22)

where "cos" is the cosine similarity between feature maps 549 (along the channel dimension), and "mean" is the mean over 550 all spatial locations. Similarly, for CLS, the warping feature 551 loss  $\ell_W$  is replaced by the similarity loss of the projected 552 features  $\ell_s$ . Note that the similarity loss  $\ell_s$  is applied after the 553 projection-thus the warping function only needs to predict 554 the residual motion in the camera view, which is the object 555 motion in time, so as to align the projected features. 556

## IV. EXPERIMENTS

We validate the effectiveness of the proposed view synchronization model on two unsynchronized multiview tasks: multiview crowd counting and multiview 3-D human pose estimation.

## 562 A. Implementation Details

The synchronization model consists of two parts: motion 563 estimation network and feature warping layer. The input of the 564 motion estimation network is the unsynchronized multiview 565 features (the concatenation of the projected features) for SLS 566 or the matching result of the 2-D camera-view features for 567 CLS, and the output is a two-channel motion flow. The layer 568 setting of the motion estimation network is shown in Table I. 569 The feature warping layer warps the features from other views 570 to align with the reference views, guided by the estimated 571 motion flow. The feature warping layer is based on the image 572 573 resampler from the spatial transformation layer in [52].

The synchronized multiview model consists of feature 574 extraction module, projection module, and multiview pre-575 diction module. For the multiview counting model [9], 576 Table II shows the model setting of the feature extraction 577 and multiview prediction module. For the 3-D pose estima-578 tion model [2], the feature extraction module consists of a 579 ResNet-152 network, a series of transposed convolution layers 580 and a  $1 \times 1$  convolution layer to predict joint heatmaps [53], 581 and the V2V-PoseNet [54] is used for multiview prediction, 582 which is based on hour-glass network [55]. 583

#### TABLE I

LAYER SETTINGS FOR THE MOTION ESTIMATION NET IN THE VIEW SYNCHRONIZATION MODULE. THE FILTER DIMENSIONS ARE OUTPUT CHANNELS, INPUT CHANNELS AND FILTER SIZE  $w_0 \times h_0$ 

| Layer  | Filter                                |
|--------|---------------------------------------|
| conv 1 | $128 \times n \times 5 \times 5$      |
| conv 2 | $128 \times 128 \times 5 \times 5$    |
| conv 3 | $64 \times 128 \times 5 \times 5$     |
| conv 4 | $64 \times 64 \times 5 \times 5$      |
| conv 5 | $32 \times 64 \times 5 \times 5$      |
| conv 6 | $2\!\times\!32\!\times\!5\!\times\!5$ |

#### TABLE II

MODEL SETTING OF THE SYNCHRONIZED MULTIVIEW COUNTING MODEL [9], CONSISTING OF FEATURE EXTRACTION AND MULTIVIEW PREDICTION. THE FILTER DIMENSIONS ARE OUTPUT CHANNELS, INPUT CHANNELS, AND FILTER SIZE  $(w \times h)$ 

| Featu   | Feature extraction               |  |  |  |  |  |
|---------|----------------------------------|--|--|--|--|--|
| 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  | $1 \times 32 \times 5 \times 5$  |  |  |  |  |  |

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

## B. Experiment Setup

We test four versions of our synchronization model: scenelevel synchronization (denoted as SLS), and CLS using concatenation, correlation, or correlation with epipolar-guided weights (denoted as CLS-cat, CLS-cor, CLS-epi) for the matching module (Section III-B.2). The synchronization models are trained with the multiview DNNs introduced in each application later.

We consider two training scenarios: 1) both synchronized 592 and unsynchronized training data is available and 2) only 593 unsynchronized training data is available, which is the more 594 difficult setting. For the first training scenario, we compare 595 against two comparison methods: BaseS trains the DNN only 596 on the synchronized data; BaseSU fine-tunes the BaseS model 597 using the unsynchronized training data (using the full training 598 set). For the second training scenario, BaseU trains the DNN 599 directly from the unsynchronized data. Note that traditional 600 synchronization methods [29]-[33] are based on videos (tem-601 poral information) and assume high-fps cameras with fixed 602 frame rates, which are unavailable in our problem setting. 603 Thus, traditional and video-based synchronization methods are 604 not suitable for comparison. 605

To test the proposed method, we first create an unsyn-606 chonized multiview dataset from the existing multiview 607 datasets (the specific datasets are introduced in each appli-608 cation later). In particular, suppose the frame sequence in the 609 reference view is captured at times  $t_0 + k\Delta t$ , where  $\Delta t$  is the 610 time offset between neighbor frames,  $k \in \{0, \dots, N-1\}$  and 611 N is the number of frames. For view i, the unsynchronized 612 frames are captured at times  $t_0 + k\Delta t + \delta_{i,k}$ , where  $\delta_{i,k}$  is 613 the desynchronization time offset between view i and the 614 reference view. We consider two settings of the desynchro-615 nization offset. The first is a *constant latency* for each view, 616  $\delta_{i,k} = \tau_i$ , for some constant value  $\tau_i$ . The second is *random* 617

584

585

586

587

588

589

590

TABLE III UNSYNCHRONIZED MULTIVIEW COUNTING: EXPERIMENT RESULTS FOR TRAINING SET WITH BOTH SYNCHRONIZED AND UNSYNCHRONIZED FRAMES. TWO DESYNCHRONIZATION SETTINGS ARE TESTED: CONSTANT LATENCY AND RANDOM LATENCY. THE EVALUA-TION METRIC IS MAE AND NAE

|                  |         | PETS2009   |                          | CityS                    | Street     |
|------------------|---------|------------|--------------------------|--------------------------|------------|
|                  |         | constant   | random                   | constant                 | random     |
| loss             | model   | MAE NAE    | MAE NAE                  | MAE NAE                  | MAE NAE    |
| 0                | BaseS   | 7.21 0.200 | 4.58 0.139               | 9.07 0.108               | 8.86 0.107 |
| $\ell_p$         | BaseSU  | 4.36 0.137 | 4.30 0.140               | 9.02 0.106               | 8.82 0.108 |
| $\ell_p, \ell_W$ | SLS     | 4.49 0.145 | 4.91 0.154               | 8.23 0.102               | 8.02 0.101 |
|                  | CLS-cat | 4.18 0.130 | 4.85 0.150               | 8.82 0.111               | 8.57 0.108 |
|                  | CLS-cor | 4.13 0.135 | 4.03 0.128               | 8.03 0.099               | 7.99 0.098 |
|                  | CLS-epi | 3.95 0.130 | <u>4.09</u> <u>0.129</u> | <u>8.05</u> <u>0.100</u> | 7.93 0.096 |

#### TABLE IV

UNSYNCHRONIZED MULTIVIEW COUNTING: EXPERIMENT RESULTS FOR TRAINING SET WITH ONLY UNSYNCHRONIZED FRAMES UNDER CON-STANT AND RANDOM LATENCY

|                  |         | PETS2009                 |                          | CityStreet               |                          |
|------------------|---------|--------------------------|--------------------------|--------------------------|--------------------------|
|                  |         | constant                 | random                   | constant                 | random                   |
| loss             | model   | MAE NAE                  | MAE NAE                  | MAE NAE                  | MAE NAE                  |
| $\ell_p$         | BaseU   | 6.18 0.187               | 6.22 0.192               | 10.22 0.134              | 9.35 0.121               |
| $\ell_p, \ell_s$ | SLS     | 5.37 0.178               | 4.82 0.150               | 8.50 0.105               | 8.33 0.100               |
|                  | CLS-cat | 6.00 0.186               | 6.08 0.189               | 8.48 0.102               | 9.17 0.110               |
|                  | CLS-cor | <b>4.18</b> <u>0.136</u> | 4.34 0.136               | <b>8.02</b> <u>0.098</u> | <u>7.77</u> 0.093        |
|                  | CLS-epi | <u>4.25</u> 0.135        | <u>4.77</u> <u>0.144</u> | <u>8.04</u> 0.095        | <b>7.70</b> <u>0.094</u> |
| $\ell_p$         | SLS     | 7.13 0.226               | 5.30 0.162               | 8.77 0.107               | 8.45 0.107               |
|                  | CLS-cat | 6.30 0.194               | 5.98 0.184               | 8.28 <u>0.098</u>        | 9.15 0.108               |
|                  | CLS-cor | <b>4.25</b> 0.138        | 4.49 0.141               | <u>8.20</u> 0.099        | <u>8.10</u> 0.102        |
|                  | CLS-epi | 4.27 <b>0.135</b>        | 4.53 0.143               | 8.16 0.097               | 7.86 0.096               |

<sup>618</sup> *latency*, where the offset for each frame and view is randomly <sup>619</sup> sampled from a uniform distribution,  $\delta_{i,k} \sim U(-\kappa_i, \kappa_i)$ . <sup>620</sup> Finally, since the synchronization is with the reference view, <sup>621</sup> the ground-truth labels for the multiview task correspond to <sup>622</sup> the times of the reference view,  $t_0 + k\Delta t$ .

## 623 C. Unsynchronized Multiview Counting

We first apply our synchronization model to unsynchronized 624 multiview counting system, whose bandwidth is assumed to be 625 limited and the frame latency between cameras can be fixed 626 or random. Here, we adopt the multiview multiscale fusion 627 model (MVMS) from [9], which is the state-of-the-art model 628 for multiview counting DNNs. We embed the synchronization 629 models in the MVMS model to handle the unsynchronized 630 multiview frames for crowd counting. 631

*Datasets and Metric:* Two multiview counting datasets
 used in [9], PETS2009 [56] and CityStreet [9], are selected
 and desynchronized for the experiments.

PETS2009 contains three views (cameras 1, 2, and 3), and 635 the first camera view is chosen as the reference view. The 636 input image resolution  $(w \times h)$  is 384  $\times$  288 and the ground-637 truth scene-level density map resolution is  $152 \times 177$ . There 638 are 825 multiview frames for training and 514 frames for 639 testing. The frame rate of PETS2009 is 7 fps ( $\Delta t = 1/7s$ ). For 640 constant frame latency,  $\tau_i \in \{5 \text{ s}, -5 \text{ s}\}$  is used for cameras 641 2 and 3, and  $\kappa_i = 5$  s for random latency. 642

*CityStreet* proposed in [9] consists of three views (cameras 1, 3, and 4), and camera 1 is chosen as the reference view. The input image resolution is  $676 \times 380$  and the ground-truth density map resolution is  $160 \times 192$ . There are 500 multiview frames, and the first 300 are used for training and the

#### TABLE V

UNSYNCHRONIZED MULTIVIEW COUNTING: EXPERIMENT RESULTS FOR TRAINING SET WITH ONLY UNSYNCHRONIZED FRAMES UNDER CON-STANT AND RANDOM LATENCY AND USING GROUND-TRUTH CAL-CULATED FROM UNSYNCHRONIZED MULTIVIEW FRAMES

|                          |         | PETS2009    |             | CityStreet  |             |
|--------------------------|---------|-------------|-------------|-------------|-------------|
|                          |         | constant    | random      | constant    | random      |
| loss                     | model   | MAE NAE     | MAE NAE     | MAE NAE     | MAE NAE     |
| $\ell_p$                 | BaseU   | 14.89 0.458 | 10.95 0.484 | 10.96 0.146 | 11.30 0.149 |
|                          | SLS     | 6.80 0.229  | 6.58 0.283  | 9.18 0.111  | 9.49 0.117  |
|                          | CLS-cat | 7.41 0.237  | 6.10 0.237  | 9.72 0.130  | 9.69 0.129  |
| $\epsilon_p, \epsilon_s$ | CLS-cor | 5.91 0.201  | 5.93 0.240  | 8.55 0.106  | 8.31 0.107  |
|                          | CLS-epi | 5.72 0.184  | 4.80 0.187  | 8.32 0.104  | 8.05 0.102  |
|                          | SLS     | 7.85 0.274  | 7.22 0.313  | 9.31 0.109  | 8.91 0.108  |
| $\ell_p$                 | CLS-cat | 7.52 0.240  | 6.20 0.243  | 8.48 0.107  | 9.85 0.121  |
|                          | CLS-cor | 6.98 0.244  | 6.26 0.282  | 8.03 0.099  | 8.24 0.107  |
|                          | CLS-epi | 6.80 0.229  | 5.18 0.200  | 8.23 0.102  | 8.16 0.103  |

remaining 200 for testing. The frame rate of CityStreet is 1 fps ( $\Delta t = 1$  s).<sup>2</sup> For constant latency,  $\tau_i \in \{3 \text{ s}, -3 \text{ s}\}$  for cameras 3 and 4, and  $\kappa_i = 3$  s for random latency. (649)

Following [9], the mean absolute error (MAE) and normalized absolute error (NAE) of the predicted counts on the test set are used as the evaluation metric 653

$$MAE = \frac{1}{N} \sum_{i=1}^{N} |\hat{c}_i - c_i|$$
(23) 65

NAE = 
$$\frac{1}{N} \sum_{i=1}^{N} |\hat{c}_i - c_i| / c_i$$
 (24) 65

where  $c_i$  is the ground truth count and  $\hat{c}_i$  is the predicted count, and N is the number of testing images.

2) Results for Training With Synced and Unsynced Frames: 658 The experimental results using training with synchronized 659 and unsynchronized frames are shown in Table III. The 660 hyperparameter  $\gamma = 1$  is used for feature warping loss. 661 On both datasets, our CLS methods, CLS-cor and CLS-epi, 662 perform better than other methods, including the baselines, 663 demonstrating the efficicacy of our approach. SLS performs 664 worse than CLS methods, due to the ambiguity of the projected 665 features from multiviews. Furthermore after projection to the 666 ground-plane, the crowd movement between frames  $I_i^{t_0}$  and  $I_i^{t_i}$ 667 on the ground-plane is less salient due to the low resolution of 668 the ground-plane feature map. CLS-cat performs worse among 669 the CLS methods because simple concatenation of features 670 cannot capture the image correspondence between different 671 views to estimate the motion flow. Finally, the two baselines 672 (BaseS and BaseSU) perform badly on CityStreet because 673 of the larger scene with larger crowd movement between 674 neighboring frames (due to lower frame rate). 675

3) Results for Training With Only Unsynchronized Frames: 676 The experiment results by training with only unsynchronized 677 frames (which is a more practical real-world case) are shown 678 in Table IV. Since the synchronized frames are not available, 679 the MVMS model weights are trained from scratch using only 680 unsynchronized data. Our models are trained with the similar-681 ity loss  $\ell_s$  (with hyperparameter  $\gamma = 1000$ ), which encourages 682 alignment of the projected multiview features. Generally, with-683 out the synchronized frames in the training stage, the counting 684

<sup>2</sup>We obtained the higher fps version from the dataset authors.

#### TABLE VI

Ablation Study on the Multiscale Architecture of the Proposed Methods for Multiview Counting on CityStreet Dataset. The Top Rows Show Performance When Training With Synchronized and Unsynchronized Frames and Using Feature Warping Loss  $\ell_W$ . The Bottom Is Training Only on Unsynchronized Frames Using Feature Similarity Loss  $\ell_s$ 

| Loss/Training data         | Method  | Multi | -scale | Single-scale |       |
|----------------------------|---------|-------|--------|--------------|-------|
|                            |         | MAE   | NAE    | MAE          | NAE   |
|                            | SLS     | 8.02  | 0.101  | 8.31         | 0.100 |
| $\ell_p, \ell_W$ /         | CLS-cat | 8.57  | 0.108  | 8.77         | 0.102 |
| sync and unsync            | CLS-cor | 7.99  | 0.098  | 8.25         | 0.099 |
|                            | CLS-epi | 7.93  | 0.096  | 8.12         | 0.098 |
|                            | SLS     | 8.33  | 0.100  | 8.95         | 0.112 |
| $\ell_p, \ell_s$<br>unsync | CLS-cat | 9.17  | 0.110  | 9.54         | 0.116 |
|                            | CLS-cor | 7.77  | 0.093  | 8.62         | 0.111 |
|                            | CLS-epi | 7.70  | 0.094  | 8.59         | 0.110 |

#### TABLE VII

AVERAGE FEATURE MAPS VALUE MEAN AND VARIANCE BEFORE AND AFTER THE FEATURE WARPING OF VIEWS 2 AND 3 OF CITYSTREET DATASET

| Method         | view 2            | view 3            |
|----------------|-------------------|-------------------|
| before warping | $0.686 \pm 1.006$ | $0.777 \pm 0.704$ |
| after warping  | $0.670\pm0.976$   | $0.761 \pm 0.692$ |

error increases for each method. Nonetheless, the proposed 685 CLS models CLS-cor and CLS-epi perform much better than 686 the baseline BaseU. CLS-cor and CLS-epi trained on only 687 unsynchronized data also performs better (on CityStreet) or on 688 par with (on PETS2009) the baseline BaseSU, which uses both 689 synchronized and unsynchronized training data. These two 690 results demonstrate the efficacy of our synchronization model 691 when only unsynchronized training data are available. Finally, 692 the error for almost all synchronization models increases on 693 both datasets when training without the similarity loss ( $\ell_p$  in 694 Table IV). This demonstrates the effectiveness of using  $\ell_s$  to 695 align the multiview features in training. 696

4) Results for Using Ground-Truth From Unsynchronized 697 Multiview Images: In the previous experiments (training 698 with only unsynchronized frames, see Section IV-C.3), the 699 ground-truth is corresponded (synchronized) to the frames 700 of the reference view. We also perform experiments when 701 the ground-truth scene-level density maps are calculated 702 from the unsynchronized multiview images. Specifically, 703 we project the same person's image coordinates of each 704 unsynchronized view to the world plane and the average of the 705 projection results is used as the ground-truth person location 706 on the ground. Then, we use the obtained person location map 707 to generate the scene-level density map. 708

The results for training with ground-truth from unsynchronized multiview images and only unsynchronized frames can be seen in Table V. From the table, we can also find that the proposed method CLS-cor/CLS-epi can achieve better performance than other methods and CLS-epi achieves the best performance, and the performance can be further improved by adding similarity loss  $\ell_s$ .

5) Ablation Study on the Multiscale Architecture: We next
present an ablation study on the multiscale architecture for
the multiview counting in Table VI. Generally, the multiscale architecture performs better than single-scale architecture
models, and the proposed CLS-cor/CLS-epi can perform better

#### TABLE VIII

COMPARISON OF METHODS ON THE CITYSTREET DATASET WITH ONLY UNSYNCHRONIZED FRAMES (BOTH CONSTANT AND RANDOM UNSYNCHRONIZED FRAMES)

| Method           | constant           | random             |
|------------------|--------------------|--------------------|
| Color correction | 8.90/0.108         | 8.64/0.100         |
| SLS              | 8.50/0.105         | 8.33/0.100         |
| CLS-cat          | 8.48/0.102         | 9.17/0.110         |
| CLS-cor          | <b>8.02</b> /0.098 | 7.77/ <b>0.093</b> |
| CLS-epi          | 8.04/ <b>0.095</b> | <b>7.70</b> /0.094 |

## TABLE IX

MODEL PARAMETER NUMBER AND RUNNING SPEED COMPARISON OF THE BASELINE METHODS BASES/BASESU/BASEU AND THE PROPOSED SLS, CLS-CAT, CLS-COR, AND CLS-EPI FOR MULTIVIEW COUNT-ING ON CITYSTREET DATASET. THE INPUT RESOLUTION FOR THE CORRELATION STEP OF THE CAMERA-VIEW SYNCHRO-NIZATION MODULE IS 160 × 95

| Method             | Paras. Num | FPS  |
|--------------------|------------|------|
| BaseS/BaseSU/BaseU | 853.4K     | 21.9 |
| SLS                | 3.7M       | 8.3  |
| CLS-cat            | 3.7M       | 8.9  |
| CLS-cor            | 37.3M      | 7.2  |
| CLS-epi            | 37.3M      | 3.6  |

than SLS or CLS-cat under both single-scale or multiscale 721 architecture, and under both training paradigms (sync and 722 unsynced, or only unsynced). 723

6) Ablation Study on Color Correlation: The feature warp-724 ing module only applies spatial shifting on the features of 725 the unsynced views, i.e., it does not change the values 726 (e.g., color) of the unsynced features [see (5) and (10)]. 727 To demonstrate this, we calculate the average statistics (mean 728 and variance) of the feature maps before and after feature 729 warping of Views 2 and 3 of CityStreet, and present the 730 results in Table VII. The statistics of the feature maps do not 731 change much after performing feature warping, and thus the 732 performance improvement of the feature warping module is 733 not due to color correction (feature value changes). 734

We further perform an ablation study to show that image 735 color correction by itself cannot solve the frame desynchro-736 nization problem. On the CityStreet dataset, in the baseline 737 model (MVMS [9]), we add a learnable "color correction" 738 layer, comprising an extra  $1 \times 1$  convolution layer (32) 739 channels) in the branches of the other camera views before 740 the projection and fusion step. The results are denoted as 741 "color correction" in Table VIII. The error for using "color 742 correction" is worse than the proposed SLS, CLS-cor, and 743 CLS-epi. The reason is that the desynchronization issue comes 744 from the capture time difference between camera views, which 745 is better solved by spatial shifting of features rather than color 746 correction (changing feature values). 747

7) Model Size and Running Speed Comparison: We present 748 the model size (number of parameters) and running speed 749 of the baseline methods and the proposed SLS, CLS-cat, 750 CLS-cor, and CLS-epi in Table IX. The input resolution 751 for the correlation step of the camera-view synchronization 752 module is  $160 \times 95$ . All models are tested on the CityStreet 753 dataset with a NVIDIA 1080Ti GPU. The baseline methods 754 (BaseS, BaseSU, and BaseU) do not use view synchronization 755 modules, so their model sizes are smaller and running speeds 756 are faster. The proposed CLS-cor and CLS-epi methods have 757

TABLE XUNSYNCHRONIZED 3-D HUMAN POSE ESTIMATION: EXPERIMENTRESULTS WITH RANDOM LATENCY. FOR "CLS-COR" AND "CLS-EPI,"THE CONSISTENCY LOSS HYPERPARAMETER  $\gamma = 0.01$ . THE EVAL-<br/>UATION METRIC IS MPJPE AND ABSOLUTE<br/>POSITION MPJPE (LEFT/RIGHT)

| Latency             | 8/50s     | 32/50s    | 64/50s      |
|---------------------|-----------|-----------|-------------|
| BaseS               | 62.8/59.2 | 78.6/78.2 | 151.1/151.5 |
| BaseSU              | 26.5/27.8 | 49.9/50.1 | 69.4/69.2   |
| BaseU               | 37.3/38.9 | 50.9/50.6 | 71.0/70.7   |
| $CLS-cor(\gamma=0)$ | 25.8/26.9 | 36.5/36.7 | 56.6/56.9   |
| CLS-cor             | 25.8/27.0 | 38.2/38.7 | 46.8/47.1   |
| CLS-epi             | 25.7/26.8 | 37.6/37.8 | 45.7/45.6   |

## TABLE XI

Detailed Performance for Unsynchronized 3-D Human Pose Estimation With Random Latency  $\kappa_i = 8/50$  s. The Evaluation Metric Is MPJPE

| Pose            | BaseS | BaseSU | BaseU | CLS-cor( $\gamma$ =0) | CLS-cor | CLS-epi |
|-----------------|-------|--------|-------|-----------------------|---------|---------|
| Directions      | 42.8  | 29.3   | 34.3  | 26.1                  | 25.8    | 26.1    |
| Discussion      | 60.7  | 28.4   | 38.8  | 27.3                  | 26.7    | 27.0    |
| Eating          | 60.7  | 26.4   | 28.8  | 23.9                  | 24.0    | 23.4    |
| Greeting        | 63.8  | 19.7   | 32.3  | 25.3                  | 24.3    | 25.1    |
| PhoneCall       | 52.2  | 25.7   | 31.0  | 24.7                  | 24.5    | 24.4    |
| Posing          | 49.7  | 22.0   | 27.6  | 24.1                  | 24.0    | 24.0    |
| Purchases       | 67.5  | 24.4   | 52.5  | 28.7                  | 27.4    | 28.8    |
| Sitting         | 33.2  | 22.6   | 36.6  | 23.8                  | 24.0    | 24.0    |
| SittingDown     | 37.4  | 25.7   | 66.6  | 25.9                  | 26.8    | 27.2    |
| Smoking         | 42.2  | 25.7   | 31.2  | 24.8                  | 24.3    | 24.4    |
| TakingPhoto     | 59.9  | 24.3   | 44.2  | 28.2                  | 27.9    | 27.2    |
| Waiting         | 44.3  | 19.5   | 35.8  | 23.2                  | 23.8    | 24.2    |
| Walking         | 161.1 | 31.9   | 32.1  | 27.0                  | 30.2    | 27.8    |
| WalkingDogs     | 91.5  | 34.2   | 54.8  | 30.1                  | 30.1    | 29.8    |
| WalkingTogether | 126.8 | 33.9   | 31.8  | 25.5                  | 26.8    | 25.5    |
| Average         | 62.8  | 26.5   | 37.3  | 25.8                  | 25.8    | 25.7    |

the correlation module, and thus have more parameters than
SLS or CLS-cat. CLS-epi is slower than CLS-cor due to the
extra multiplication step with the epipolar weights.

8) Visualization Results: Example results are shown in 761 Fig. 5. Generally, the proposed synchronization methods 762 CLS-epi and CLS-cor can predict better quality density maps, 763 such as in the red box regions in the figure, where comparison 764 methods tend to over-count these regions due to the same per-765 son being counted multiple times in unsynchronized frames. 766 Furthermore, we also observe that the predicted density map 767 is with better quality when synchronized frames are avail-768 able compared to training with only unsynchronized frames. 769 Finally, the prediction results are improved if similarity loss 770 is enforced when training with only unsynchronized frames, 771 such as the methods CLS-epi and CLS-cor on PETS2009. 772

## 773 D. Unsynchronized 3-D Human Pose Estimation

774 We next apply our synchronization model to the unsynchronized 3-D human pose estimation task. The DNNs model 775 for 3-D human pose estimation is adopted from [2], which 776 proposed two learnable triangulation methods for multiview 777 3-D human pose from multiple 2-D views: algebraic triangu-778 lation and volumetric aggregation. Here, we use volumetric 779 aggregation (softmax aggregation) as the multiview fusion 780 DNN in the experiments. 781

*Datasets and Metrics:* We use the Human3.6M [57]
dataset, which consists of 3.6 million frames from four synchronized 50 Hz digital cameras along with the 3-D pose

TABLE XII

Detailed Performance for Unsynchronized 3-D Human Pose Estimation With Random Latency  $\kappa_i = 32/50$  s. The Evaluation Metric Is MPJPE

| Pose            | BaseS | BaseSU | BaseU | CLS-cor( $\gamma=0$ ) | CLS-cor | CLS-epi |
|-----------------|-------|--------|-------|-----------------------|---------|---------|
| Directions      | 46.2  | 48.7   | 65.1  | 42.9                  | 42.5    | 43.9    |
| Discussion      | 75.6  | 53.6   | 55.2  | 38.9                  | 41.0    | 41.6    |
| Eating          | 64.5  | 39.1   | 40.2  | 32.5                  | 32.7    | 30.8    |
| Greeting        | 71.5  | 48.5   | 55.4  | 35.7                  | 38.1    | 36.8    |
| PhoneCall       | 64.5  | 43.6   | 43.0  | 33.9                  | 35.2    | 35.1    |
| Posing          | 49.3  | 42.1   | 43.3  | 32.7                  | 33.3    | 30.8    |
| Purchases       | 111.5 | 50.9   | 48.7  | 35.9                  | 42.4    | 40.4    |
| Sitting         | 55.2  | 46.0   | 43.7  | 33.6                  | 33.8    | 34.7    |
| SittingDown     | 108.3 | 79.3   | 64.9  | 36.8                  | 41.8    | 42.8    |
| Smoking         | 54.5  | 44.3   | 44.0  | 35.5                  | 35.9    | 35.7    |
| TakingPhoto     | 87.9  | 57.0   | 58.6  | 39.3                  | 43.0    | 41.3    |
| Waiting         | 64.3  | 45.6   | 47.3  | 35.5                  | 33.7    | 35.0    |
| Walking         | 150.6 | 47.6   | 48.1  | 34.2                  | 37.1    | 34.2    |
| WalkingDogs     | 123.1 | 66.2   | 67.5  | 44.5                  | 49.2    | 49.1    |
| WalkingTogether | 125.5 | 50.3   | 52.7  | 36.9                  | 38.5    | 34.9    |
| Average         | 78.6  | 49.9   | 50.9  | 36.5                  | 38.2    | 37.6    |

#### TABLE XIII

Detailed Performance for Unsynchronized 3-D Human Pose Estimation With Random Latency  $\kappa_i = 64/50$  s. The Evaluation Metric Is MPJPE

| Pose            | BaseS | BaseSU | BaseU | CLS-cor( $\gamma=0$ ) | ) CLS-cor | CLS-epi |
|-----------------|-------|--------|-------|-----------------------|-----------|---------|
| Directions      | 99.2  | 83.2   | 76.3  | 70.3                  | 64.8      | 66.5    |
| Discussion      | 144.1 | 72.0   | 67.5  | 57.3                  | 48.2      | 48.4    |
| Eating          | 138.2 | 55.3   | 63.2  | 44.7                  | 40.4      | 37.9    |
| Greeting        | 181.3 | 68.0   | 74.1  | 54.8                  | 46.9      | 46.3    |
| PhoneCall       | 138.1 | 58.8   | 61.1  | 49.2                  | 40.7      | 40.5    |
| Posing          | 121.7 | 53.5   | 50.2  | 42.1                  | 36.6      | 36.3    |
| Purchases       | 155.7 | 69.0   | 62.4  | 58.1                  | 47.0      | 50.6    |
| Sitting         | 74.0  | 64.2   | 67.8  | 55.2                  | 41.1      | 39.6    |
| SittingDown     | 103.8 | 112.3  | 140.7 | 89.8                  | 54.6      | 50.6    |
| Smoking         | 112.7 | 58.7   | 60.3  | 49.2                  | 41.7      | 41.7    |
| TakingPhoto     | 166.8 | 76.8   | 79.6  | 64.5                  | 57.7      | 53.5    |
| Waiting         | 120.2 | 62.7   | 61.1  | 51.1                  | 42.0      | 42.6    |
| Walking         | 301.1 | 66.3   | 69.2  | 49.7                  | 44.0      | 41.7    |
| WalkingDogs     | 219.2 | 95.9   | 91.6  | 77.0                  | 62.7      | 55.5    |
| WalkingTogether | 302.7 | 67.2   | 68.9  | 54.5                  | 43.5      | 42.5    |
| Average         | 151.1 | 69.4   | 71.0  | 56.6                  | 46.8      | 45.7    |

#### TABLE XIV

UNSYNCHRONIZED 3-D HUMAN POSE ESTIMATION: CLS-EPI EXPERIMENT RESULTS WITH DIFFERENT HYPERPARAMETER  $\gamma$ . THE EVALUATION METRIC IS MPJPE

| $\gamma$            | 0.005 | 0.01 | 0.02 |
|---------------------|-------|------|------|
| $\kappa_i = 8/50s$  | 25.6  | 25.7 | 26.0 |
| $\kappa_i = 32/50s$ | 38.3  | 37.6 | 37.9 |
| $\kappa_i = 64/50s$ | 51.7  | 45.7 | 46.8 |

annotations. We follow the preprocessing step<sup>3</sup> recommended 785 in [57], and sample one of every 64 frames ( $\Delta t = 64/50$ ) 786 for the testing set, and sample one of every four frames 787  $(\Delta t = 4/50)$  as the training set. The first camera view is 788 always used as the reference view (if the first camera view is 789 missing, the second one is used). We test desynchronization 790 via random frame latency, with  $\kappa_i \in \{8/50, 32/50, 64/50\}$  s, 791 and only use unsynchronized data for training. Following [2], 792 mean per point position error (MPJPE) and absolute position 793 MPJPE are used as the metric for evaluation. In training, 794 the single-view backbone uses the pretrained weights from 795 the original 3-D pose estimation model. Baseline methods 796

<sup>3</sup>https://github.com/anibali/h36m-fetch. Accessed: October 10, 2019.



Fig. 5. Examples of unsynchronized multiview crowd counting on PETS2009 (top) and CityStreet (bottom). The left shows the input multiview frames, and note that the synchronized frames (in dotted box) are not used when training with only unsynchronized frames. The input unsynchronized frames are randomly selected around the synchronized frames. For each dataset, the result of training with synchronized and unsynchronized frames  $(l_p \text{ and } l_W)$  is in row 1, the result of training only with unsynchronized frames  $(l_p \text{ and } l_W)$  is in row 1, the result of training only with unsynchronized frames  $(l_p \text{ and } l_S)$  is shown in row 2, and the result of training with unsynchronized frames and using similarity loss between projected features  $(l_p \text{ and } l_S)$  is shown in row 3. Generally, 1) the proposed synchronization methods CLS-epi and CLS-cor predict density maps with better quality compared to other comparison methods; 2) the methods achieve better performance when synchronized frames are available in training; and 3) when training only with unsynchronized frames, enforcing the similarity loss  $l_s$  can help improve the performance.

<sup>797</sup> BaseS, BaseSU and BaseU are compared with our proposed
 <sup>798</sup> camera-view synchronization models CLS-cor and CLS-epi.

2) Experiment Results: The experiments results are pre-799 sented in Table X. The original 3-D pose estimation method 800 (BaseS, BaseSU, and BaseU) cannot perform well under 801 the unsynchronized test condition, especially under large 802 latencies (e.g., 64/50 s). Our camera-view synchronization 803 methods performs better than the baseline methods, with the 804 805 performance gap increasing as the latency increases. Using similarity loss  $\ell_s$  improves the performance of our models, 806 and adding epipolar-guided weights can suppress false matches 807 and further reduces the error. The detailed performance for 808 each pose type under different frame latency settings is shown 809 in Tables XI-XIII. From the tables, we can find that the 810 proposed methods can perform especially better on the poses 811 with larger movement between unsynchronized frames, e.g., 812 Walking, WalkingDogs and WalkingTogether. 813

<sup>814</sup> 3) Ablation Study on  $\gamma$  for 3-D Pose Estimation: The <sup>815</sup> ablation study on hyperparameter  $\gamma$  for the method CLS-epi for 3-D pose estimation is presented in Table XIV. In general,  $\gamma = 0.01$  achieves better performance than other weights.

4) Model Size and Running Speed Comparison: We present 818 the model sizes and running speed comparisons of our pro-819 posed models and the baselines for 3-D pose estimation 820 in Table XV. The input resolution for the correlation step 821 of the camera-view synchronization module is  $48 \times 48$ . 822 As the original synchronized 3-D pose estimation model [2] is 823 already very large, the running speed of the proposed models 824 CLS-cor and CLS-epi is similar to the baseline methods 825 BaseS/BaseSU/BaseU. 826

5) Visualization Results: Visualization results of unsynchro-827 nized 3-D pose estimation are presented in Figs. 6 and 7. In the 828 figures, the first row shows the input unsynchronized multi-829 view frames, and the top labels indicate the unsynchronized 830 frame latency. Rows 2-8 show the 2-D key-joints projected 831 from 3-D poses of Ground-truth, BaseS, BaseSU, BaseU, 832 CLS-cor ( $\gamma = 0$ ), CLS-cor, and CLS-epi, respectively, where 833 synchronized frames are displayed for better visualization 834



Fig. 6. Examples of unsynchronized 3-D pose estimation (Walking Dogs). The first row shows the input unsynchronized multiview frames and the top labels indicate the unsynchronized frame latency (in seconds). The remaining rows show the ground-truth key joints and the predicted results. Blue lines are the 2-D key joints projected from 3-D poses, and the *synchronized* frames are used for better visualization. CLS-epi achieves the best performance among all methods, especially the prediction result of arms in view 0.



Fig. 7. Examples of unsynchronized 3-D pose estimation (Greeting). Blue lines are the 2-D key-joints projected from 3-D poses, and the synchronized frames are used for better visualization. CLS-epi achieves the best performance.

#### TABLE XV

| Method             | Paras. Num | FPS |
|--------------------|------------|-----|
| BaseS/BaseSU/BaseU | 80.6M      | 3.7 |
| CLS-cor            | 86.3M      | 3.4 |
| CLS-epi            | 86.3M      | 3.0 |

effect. In Fig. 6, BaseU fails on the unsynchronized input, and CLS-epi achieves the best performance among all methods, especially the prediction of the arms in view 1. In Fig. 7, the CLS-epi also achieves the best performance among all comparison methods.

840

## V. CONCLUSION

In this article, we focus on the issue of unsynchronized 841 cameras in DNNs-based multiview computer vision tasks. 842 We propose two view synchronization models based on single 843 frames (not videos) from each view, SLS and CLS. The two 844 models are trained and evaluated under two training settings 845 (with or without synchronized frame pairs), and a similarity 846 loss of the projected multiview features is proposed to boost 847 the performance when synchronized training pairs are not 848 available. Furthermore, to show its generality to different 849 conditions of desynchronization, the proposed models are 850 tested with desynchronization based on both constant and 851 random latency. Finally, the proposed models are applied 852 to unsynchronized multiview counting and unsynchronized 853 3-D human pose estimation, and achieve better performance 854 compared to the baseline methods. Overall, CLS model using 855 correlation and epipolar weights (CLS-epi) performs best 856 among the proposed models. 857

In addition to unsynchronized multicamera crowd counting 858 and 3-D pose estimation, the proposed method can also be 859 applied to other multicamera vision tasks, such as multi-860 camera detection [7], multicamera tracking [18]. In these 861 tasks, multicameras may also be unsynchronized due to no 862 synchronization clock or limited network bandwidth. As these 863 DNN models [7], [18] generally follow the three-stage pipeline 864 (single-view feature extraction, multiview projection and 865 fusion, and prediction), our proposed synchronization modules 866 can be inserted to adapt them to unsynchronized frames. 867

In our current model, image content matching is used 868 for view synchronization, while the 2-D-to-3-D projection 869 for multiview fusion relies on known camera parameters. 870 The multicamera surveillance tasks themselves require known 871 calibration for better multiview fusion. Note that our pro-872 posed view synchronization module based on correlation maps 873 (CLS-cor) does not require camera calibrations due to the 874 single-frame basis, and still achieves good performance. When 875 the calibrations are provided, epipolar constraint can be uti-876 lized to achieve better results (CLS-epi). In future work, the 877 2-D-3-D projection in the original multiview models could 878 be replaced with camera self-calibration modules, which can 879 allow the model to handle unsynchronized and uncalibrated 880 multicameras. 881

## References

- [1] H. Aghajan and A. Cavallaro, *Multi-Camera Networks: Principles and Applications*. New York, NY, USA: Academic, 2009.
- [2] K. Iskakov, E. Burkov, V. Lempitsky, and Y. Malkov, "Learnable triangulation of human pose," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 7718–7727.
- [3] E. Remelli, S. Han, S. Honari, P. Fua, and R. Wang, "Lightweight multiview 3D pose estimation through camera-disentangled representation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 6040–6049.
- [4] A. Kar, C. Háne, and J. Malik, "Learning a multi-view stereo machine," in *Proc. NIPS*, 2017, pp. 365–376.
- [5] P.-H. Huang, K. Matzen, J. Kopf, N. Ahuja, and J.-B. Huang, "Deep-MVS: Learning multi-view stereopsis," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 2821–2830.
- [6] P. Baque, F. Fleuret, and P. Fua, "Deep occlusion reasoning for multicamera multi-target detection," in *Proc. IEEE Int. Conf. Comput. Vis.* (*ICCV*), Oct. 2017, pp. 271–279.
- [7] T. Chavdarova *et al.*, "WILDTRACK: A multi-camera HD dataset for dense unscripted pedestrian detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 5030–5039.
- [8] L. Chen, H. Ai, R. Chen, Z. Zhuang, and S. Liu, "Cross-view tracking for multi-human 3D pose estimation at over 100 FPS," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 3279–3288.
- [9] Q. Zhang and A. B. Chan, "Wide-area crowd counting via ground-plane density maps and multi-view fusion CNNs," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 8297–8306.
- [10] M. Ye, X. Lan, Q. Leng, and J. Shen, "Cross-modality person reidentification via modality-aware collaborative ensemble learning," *IEEE Trans. Image Process.*, vol. 29, pp. 9387–9399, 2020.
- [11] K. Zhou, Y. Yang, A. Cavallaro, and T. Xiang, "Omni-scale feature learning for person re-identification," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 3702–3712.
- [12] M. Ye and P. C. Yuen, "PurifyNet: A robust person re-identification model with noisy labels," *IEEE Trans. Inf. Forensics Security*, vol. 15, pp. 2655–2666, 2020.
- [13] M. Ye, J. Shen, X. Zhang, P. C. Yuen, and S.-F. Chang, "Augmentation invariant and instance spreading feature for softmax embedding," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 2, pp. 924–939, Feb. 2022.
- [14] L. Song *et al.*, "Unsupervised domain adaptive re-identification: Theory and practice," *Pattern Recognit.*, vol. 102, Jun. 2020, Art. no. 107173.
- [15] E. Zheng, D. Ji, E. Dunn, and J.-M. Frahm, "Sparse dynamic 3D reconstruction from unsynchronized videos," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 4435–4443.
- [16] X. Zhang, B. Ozbay, M. Sznaier, and O. Camps, "Dynamics enhanced multi-camera motion segmentation from unsynchronized videos," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 4668–4676.
- [17] Q. Zhang and A. B. Chan, "3D crowd counting via multi-view fusion with 3D Gaussian kernels," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 12837–12844.
- [18] Y. He, J. Han, W. Yu, X. Hong, X. Wei, and Y. Gong, "City-scale multicamera vehicle tracking by semantic attribute parsing and cross-camera tracklet matching," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2020, pp. 576–577.
- [19] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese, "3D-R2N2: A unified approach for single and multi-view 3D object reconstruction," in *Proc. ECCV.* Cham, Switzerland: Springer, 2016, pp. 628–644.
- [20] H. Xie, H. Yao, X. Sun, S. Zhou, and S. Zhang, "Pix2 Vox: Contextaware 3D reconstruction from single and multi-view images," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 2690–2698.
- [21] M. Kocabas, S. Karagoz, and E. Akbas, "Self-supervised learning of 3D human pose using multi-view geometry," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 1077–1086.
- [22] C.-H. Chen *et al.*, "Unsupervised 3D pose estimation with geometric self-supervision," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 5714–5724.
- [23] H. Joo *et al.*, "Panoptic studio: A massively multiview system for social motion capture," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 3334–3342.
- [24] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Harvesting multiple views for marker-less 3D human pose annotations," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 6988–6997.

945

- [25] K. Takahashi, D. Mikami, M. Isogawa, and H. Kimata, "Human pose as 957 958 calibration pattern: 3D human pose estimation with multiple unsynchronized and uncalibrated cameras," in Proc. IEEE/CVF Conf. Comput. Vis. 959 Pattern Recognit. Workshops (CVPRW), Jun. 2018, pp. 1775-1782. 960
- 961 [26] C. Albl, Z. Kukelova, A. Fitzgibbon, J. Heller, M. Smid, and T. Pajdla, "On the two-view geometry of unsynchronized cameras," in Proc. IEEE 962 963 Conf. Comput. Vis. Pattern Recognit. (CVPR), Jul. 2017, pp. 4847-4856.
- [27] T. Kuo, S. Sunderrajan, and B. S. Manjunath, "Camera alignment using 964 trajectory intersections in unsynchronized videos," in Proc. IEEE Int. 965 966 Conf. Comput. Vis., Dec. 2013, pp. 1121-1128.
- [28] N. Hasler, B. Rosenhahn, T. Thormahlen, M. Wand, J. Gall, and 967 H.-P. Seidel, "Markerless motion capture with unsynchronized moving 968 cameras," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., Jun. 2009, 969 pp. 224-231. 970
- C. Dai, Y. Zheng, and X. Li, "Subframe video synchronization via 971 [29] 3D phase correlation," in Proc. Int. Conf. Image Process., Oct. 2006, 972 pp. 501-504. 973
- [30] F. Pádua, R. Carceroni, G. Santos, and K. Kutulakos, "Linear sequence-974 to-sequence alignment," IEEE Trans. Pattern Anal. Mach. Intell., vol. 32, 975 no. 2, pp. 304-320, Feb. 2010. 976
- C. Lei and Y.-H. Yang, "Tri-focal tensor-based multiple video syn-[31] 977 chronization with subframe optimization," IEEE Trans. Îmage Process., 978 979 vol. 15, no. 9, pp. 2473-2480, Sep. 2006.
- [32] P. A. Tresadern and I. D. Reid, "Video synchronization from human 980 motion using rank constraints," Comput. Vis. Image Understand., 981 vol. 113, no. 8, pp. 891-906, Aug. 2009. 982
- J. Yan and M. Pollefeys, "Video synchronization via space-time interest 983 [33] 984 point distribution," in Proc. Adv. Concepts Intell. Vis. Syst., vol. 1, 2004, pp. 12-21. 985
- [34] Y. Caspi, D. Simakov, and M. Irani, "Feature-based sequence-to-986 sequence matching," Int. J. Comput. Vis., vol. 68, no. 1, pp. 53-64, 987 2006988
- [35] C. Rao, A. Gritai, M. Shah, and T. Syeda-Mahmood, "View-invariant 989 alignment and matching of video sequences," in Proc. 9th IEEE Int. 990 Conf. Comput. Vis., Oct. 2003, pp. 939-945. 991
- [36] B. Meyer, T. Stich, and M. Pollefeys, "Subframe temporal alignment of 992 non-stationary cameras," in *Proc. Brit. Mach. Vis. Conf.*, 2008, pp. 1–10. E. Imre and A. Hilton, "Through-the-lens synchronisation for heteroge-993
- 994 [37] neous camera networks," in Proc. Brit. Mach. Vis. Conf., 2012, pp. 1-11. 995
- [38] D. Pundik and Y. Moses, "Video synchronization using temporal signals 996 from epipolar lines," in Proc. ECCV, 2010, pp. 15-28. 997
- S. N. Sinha and M. Pollefeys, "Synchronization and calibration of 998 [39] camera networks from silhouettes," in Proc. 17th Int. Conf. Pattern 999 Recognit. (ICPR), Aug. 2004, pp. 116-119. 1000
- [40] T. Gaspar, P. Oliveira, Favaro, and Paolo, "Synchronization two indepen-1001 1002 dently moving cameras without feature correspondences," in Proc. Eur. Conf. Comput. Vis., Zurich, Switerland. Springer, Sep. 2014, pp. 6-12. 1003
- I. Rocco, R. Arandjelovic, and J. Sivic, "Convolutional neural network 1004 [41] architecture for geometric matching," in Proc. IEEE Conf. Comput. Vis. 1005 1006 Pattern Recognit. (CVPR), Jul. 2017, pp. 6148-6157.
- [42] S. Phillips and K. Daniilidis, "All graphs lead to Rome: Learning 1007 geometric and cycle-consistent representations with graph convolutional 1008 networks," 2019, arXiv:1901.02078. 1009
- [43] H. Altwaijry, E. Trulls, J. Hays, P. Fua, and S. Belongie, "Learn-1010 1011 ing to match aerial images with deep attentive architectures," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2016, 1012 pp. 3539-3547. 1013
- [44] 1014 T.-W. Hui, X. Tang, and C. C. Loy, "LiteFlowNet: A lightweight convolutional neural network for optical flow estimation," in Proc. IEEE/CVF 1015 Conf. Comput. Vis. Pattern Recognit., Jun. 2018, pp. 8981-8989 1016
- 1017 [45] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "FlowNet 2.0: Evolution of optical flow estimation with deep networks," 1018 in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Jul. 2017, 1019 pp. 2462-2470. 1020
- [46] M. Bai, W. Luo, K. Kundu, and R. Urtasun, "Exploiting semantic 1021 information and deep matching for optical flow," in Proc. Eur. Conf. 1022 Comput. Vis., Cham, Switzerland: Springer, 2016, pp. 154-170. 1023
- W. Yu, K. Yang, Y. Bai, H. Yao, and Y. Rui, "DNN flow: DNN feature 1024 [47] 1025 pyramid based image matching," in Proc. Brit. Mach. Vis. Conf., 2014, 1026 pp. 1-10.

- [48] A. Dosovitskiy et al., "FlowNet: Learning optical flow with convo-1027 lutional networks," in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), 1028 Dec. 2015, pp. 2758-2766. 1029
- [49] A. Ranian and M. J. Black, "Optical flow estimation using a spatial 1030 pyramid network," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. 1031 (CVPR), Jul. 2017, pp. 4161-4170. 1032
- [50] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, "PWC-Net: CNNs for optical 1033 flow using pyramid, warping, and cost volume," in Proc. IEEE/CVF 1034 Conf. Comput. Vis. Pattern Recognit., Jun. 2018, pp. 8934-8943. 1035
- [51] H.-Y. Lai, Y.-H. Tsai, and W.-C. Chiu, "Bridging stereo match-1036 ing and optical flow via spatiotemporal correspondence," in Proc. 1037 IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2019, 1038 pp. 1890-1899. 1039
- [52] M. Jaderberg et al., "Spatial transformer networks," in Proc. Adv. Neural 1040 Inf. Process. Syst., 2015, pp. 2017-2025. 1041
- [53] B. Xiao, H. Wu, and Y. Wei, "Simple baselines for human pose 1042 estimation and tracking," in Proc. ECCV, Sep. 2018, pp. 466-481. 1043
- [54] J. Y. Chang, G. Moon, and K. M. Lee, "V2 V-PoseNet: Voxel-to-voxel 1044 prediction network for accurate 3D hand and human pose estimation 1045 from a single depth map," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern 1046 Recognit., Jun. 2018, pp. 5079-5088. 1047
- A. Newell, K. Yang, and D. Jia, "Stacked hourglass networks for [55] 1048 human pose estimation," in Proc. Eur. Conf. Comput. Vis., Oct. 2016, 1049 pp. 483-499. 1050
- J. Ferryman and A. Shahrokni, "PETS2009: Dataset and challenge," [56] 1051 in Proc. 12th IEEE Int. Workshop Perform. Eval. Tracking Surveill., 1052 Dec. 2009, pp. 1-6. 1053
- [57] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, "Human3.6M: 1054 Large scale datasets and predictive methods for 3d human sensing in 1055 natural environments," IEEE Trans. Pattern Anal. Mach. Intell., vol. 36, 1056 no. 7, pp. 1325-1339, Dec. 2013. 1057



Qi Zhang (Student Member, IEEE) received the 1058 B.Eng. degree from the Huazhong University of 1059 Science and Technology (HUST), Wuhan, China, 1060 in 2014, the M.Eng. degree from the University 1061 of Chinese Academy of Sciences (UCAS), Beijing, 1062 China, in 2017, and the Ph.D. degree in com-1063 puter science from City University of Hong Kong, 1064 Hong Kong, in 2021. 1065 He was a Post-Doctoral Researcher with City 1066

University of Hong Kong, from 2021 to 2022. 1067 He currently works as the Distinguished Associate 1068

Researcher with the Guangdong Laboratory of Artificial Intelligence and 1069 Digital Economy (SZ), Shenzhen, China. His research interests include 1070 multicamera surveillance and crowd counting. 1071



Antoni B. Chan (Senior Member, IEEE) received 1072 the B.S. and M.Eng. degrees in electrical engi-1073 neering from Cornell University, Ithaca, NY, USA, 1074 in 2000 and 2001, respectively, and the Ph.D. 1075 degree in electrical and computer engineering from the University of California, San Diego (UCSD), San Diego, CA, USA, in 2008.

He is currently a Full Professor with the Department of Computer Science, City University of Hong Kong, Hong Kong. His research interests 1081 include computer vision, machine learning, pattern 1082 1083

recognition, and music analysis.