

# Anyone here?

## Smart embedded low-resolution omnidirectional video sensor to measure room occupancy

Timothy Callemein, Kristof Van Beeck and Toon Goedemé

*EAVISE - KU Leuven*

Sint-Katelijne-Waver, BELGIUM

{firstname.lastname}@kuleuven.be

**Abstract**—In this paper, we present a room occupancy sensing solution with unique properties: (i) It is based on an omnidirectional vision camera, capturing rich scene info over a wide angle, enabling to count the number of people in a room and even their position. (ii) Although it uses a camera-input, no privacy issues arise because its extremely low image resolution, rendering people unrecognisable. (iii) The neural network inference is running entirely on a low-cost processing platform embedded in the sensor, reducing the privacy risk even further. (iv) Limited manual data annotation is needed, because of the self-training scheme we propose. Such a smart room occupancy rate sensor can be used in e.g. meeting rooms and flex-desks. Indeed, by encouraging flex-desking, the required office space can be reduced significantly. In some cases, however, a flex-desk that has been reserved remains unoccupied without an update in the reservation system. A similar problem occurs with meeting rooms, which are often under-occupied. By optimising the occupancy rate a huge reduction in costs can be achieved. Therefore, in this paper, we develop such system which determines the number of people present in office flex-desks and meeting rooms. Using an omnidirectional camera mounted in the ceiling, combined with a person detector, the company can intelligently update the reservation system based on the measured occupancy. Next to the optimisation and embedded implementation of such a self-training omnidirectional people detection algorithm, in this work we propose a novel approach that combines spatial and temporal image data, improving performance of our system on extreme low-resolution images.

**Index Terms**—privacy preserving, occupancy detection, omnidirectional, deep learning, low resolution

### I. INTRODUCTION

Companies often require larger facilities as their number of employees increases. By encouraging flex-desking, the required office space can be reduced significantly. However, the growing amount of meetings and the rise in popularity of flex-desks in many cases result in building capacity inefficiency problems: reserved meeting rooms and flex desks remain unoccupied due to people working from a different location, cancellations, rescheduling without adjusting the reservation info, and so on. In other cases, large meeting rooms get reserved, while another smaller meeting room also meets capacity requirements.

To partially resolve this issue, a passive infrared (PIR) sensor, that measures the change in reflected infrared light, can

be used to determine human activity in a room. However, the sensor proves to be inadequate since enough movement must occur for it to operate and it is unusable to specify the degree of occupancy, i.e. the number of people in the room, or even to determine which desks are taken. By placing a camera system, more information is gathered that can be analysed with greater accuracy. Omnidirectional cameras are gaining popularity in security applications because of their wide field-of-view and ease of installation. Compared to traditional cameras, they capture a 360 degree image without the need of camera re-positioning, which is the case for most traditional motorised (PTZ) cameras. While they can provide a complete overview at one glance, the images suffer from severe fish-eye lens distortion. This is not a problem when the camera images are analysed by humans. However, out-of-the-box state-of-the-art computer vision algorithms which are trained on frontal, upright persons will fail on these images. Hence, we need to retrain such a detector with similar omnidirectional image material.

Unfortunately, the available amount of omnidirectional training data with person annotations is limited. To overcome this challenge, we propose a self-training approach that uses state-of-the-art person detectors on unwarped omnidirectional data, to automatically generate annotation labels. These annotation labels are then used as training data to train a new model, capable of determining the room occupancy directly on omnidirectional images. This is illustrated in fig. 2. In our application, we want to exploit the fact that the ceiling-mounted camera is static and allow the trained models to be environment-specific.

An eternal issue with camera-based sensors is the concern of people’s privacy. Indeed, most employees do not feel comfortable when being constantly observed by cameras, and in many cases recording identifiable people in their work environment is not allowed due to legal regulations. In the system we propose, the privacy is guaranteed because of two reasons. Firstly, these privacy issues are avoided if the sensitive information is processed locally (for example on an embedded platform), and only the anonymous meta-data is outputted. Our application therefore will be optimized such that it is capable of running on an embedded platform, e.g. a Raspberry Pi. Secondly, and most importantly, our resulting system works

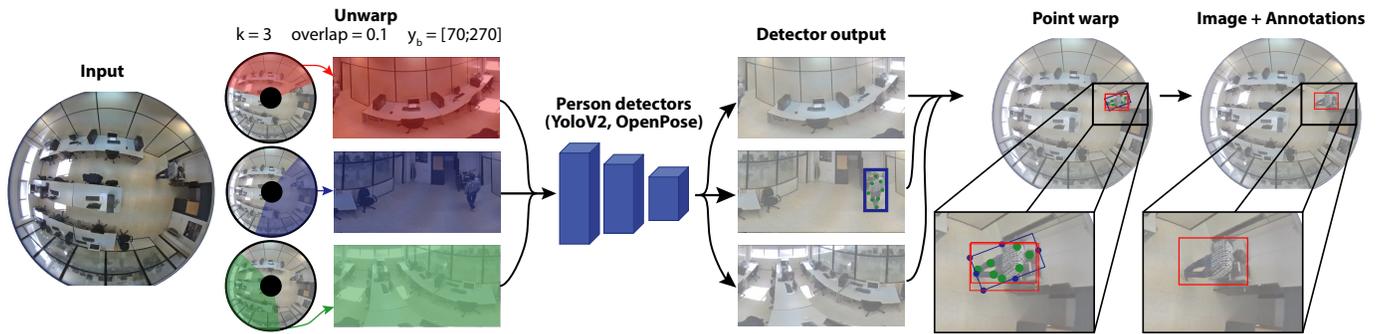


Fig. 1: Overview of our proposed self-training approach.

on extreme low image resolution data, in which people are inherently unable to be recognised. The work by Butler *et al.* [1] supports this, indicating that the sense of privacy is increased when the image contains less details (for example, by lowering the resolution). After the self-training step, we can even turn the omnidirectional lens out-of-focus, yielding identical downscaled low-resolution input images, but making hacking the sensor purposeless. Figure 3 shows example frames, the two leftmost frames showing the high and low resolution frame with a lens in focus, and the two rightmost frames with a lens places out of focus. Apart from regulations, the awareness of being recorded can be considered obtrusive and induce the feeling of being watched and monitored. Yet by designing the outer shape of the sensor to not resemble a camera, this feeling is greatly reduced.

For the above reasons, we propose to train a state-of-the-art object detector (YOLOv2 [2]) on low-resolution omnidirectional data. First, we lower the resolution and needed computational power by decreasing the network resolution. To overcome network architectural limitations, an additional effort was made to reduce the image resolution further. While lowering the resolution increases privacy, the loss in data will increase the challenge of to accurately detect the room occupancy. We therefore propose a novel approach

that incorporates temporal information to compensate the loss in spatial information. Towards this goal, we train several object detection models with several different input and image augmentation settings. In our application, where we aim to count the degree of occupancy, and therefore pay less attention to the location and bounding box output of our models, we use a count-by-detection methodology as end result.

This work goes beyond previous work by Callemain *et al.* [3], with the important novelty that the neural network is implemented on a low-cost embedded device, after several optimisations. Moreover, our combination of spatial and temporal image data is clearly boosting the detection performance and further reduce the input resolution as compared to their result.

The remainder of this paper is structured as follows. Section II) discusses the related work followed by section III describing our suggested approach to generating object detection labels and how we can further reduce the image resolution by combining temporal information. Followed by section IV where we evaluate our approach on three different datasets. In section V our conclusion with a discussion and possible future work.

## II. RELATED WORK

Object counting has many applications and challenges. Our case mainly focuses on person counting using omnidirectional camera images. Intuitively, we would look at crowd counting architectures capable of estimating the number of people in crowds [4]–[6]. Such techniques will train a convolutional neural network (CNN) model capable of estimating a crowd density map, based on the head location. These techniques have no interest in the exact spatial location of the crowd and only focus on estimating the number of people. However, most crowd counting techniques expect dense crowds on which a relative low error is allowed compared to the high value output. Our case however, only has a limited sparse number of people (12 max), which is far more sensitive to these errors.

Instead of room occupancy detection using crowd counting techniques, one might also use an object detector and simply count the number of detections. For example, work by Zou *et al.* [7] uses a two step approach (temporal histogram of gradient (HoG) head model followed by CNN head classifier) to determine the degree of occupancy. However, their approach does

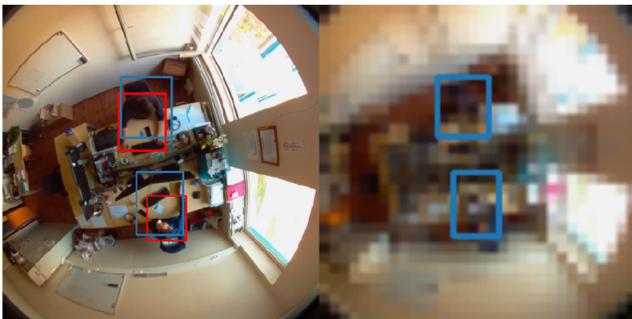
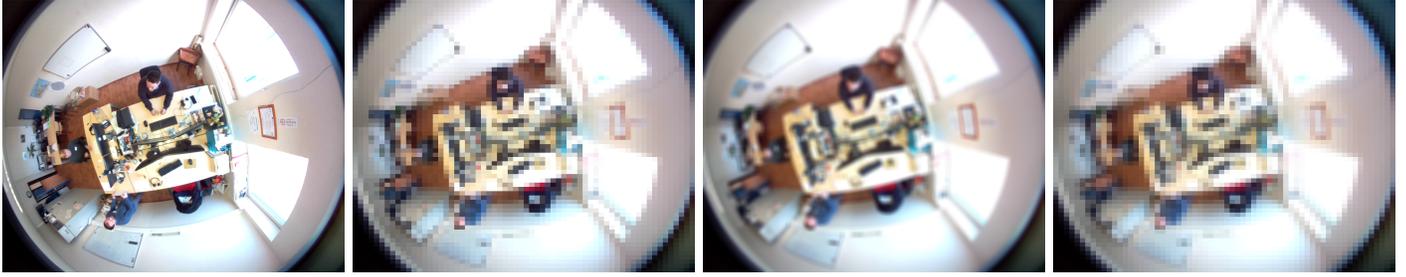


Fig. 2: Private dataset example frame, high resolution (left) and low-resolution (right), with automatically generated annotations (red) and detections based on low-resolution image (blue)



(a) High-resolution focused

(b) Low-resolution focused

(c) High-resolution out-of-focus

(d) Low-resolution out-of-focus

Fig. 3: Example camera frames, one high and low resolution frame (64 px) for a frame with a lens in focus, another with a out of focus lens

not focus on privacy concerns and uses wall mounted cameras instead which are sensitive to scene occlusions. Other work by Newsham *et al.* compares a wide variety of sensors mounted on top of a computer screen to determine the occupancy degree (including thermal, PIR and radar sensors). These sensors avoid the recording of privacy sensitive data. However, they have several disadvantages: they require installing additional (costly) hardware, enough movement is needed for reliable detection and they are unable to determine the exact degree of occupancy. In this work, we aim to remain as unobtrusive as possible, using only a single wide angle ceiling mounted camera. Therefore, to determine the occupancy degree, we propose to use a vision-based person detector. Our approach is able to detect persons in extremely low-resolution images, such that the privacy is inherently preserved.

Several deep learning architectures are capable of efficient person detection [2], [8], [9]. Often the first layers are trained on a large scale dataset (ImageNet [10]) and afterwards the full network is fine-tuned for object detection on smaller datasets (VOC [11], COCO [12]). Our application will use omnidirectional cameras producing heavily distorted images that are not included in these datasets. To overcome this challenge, Seidel *et al.* [13] proposes an approach that first transforms the omnidirectional images to perspective images. On these perspective images they use a person detector, and compare different non maximum suppression (NMS) approaches to combine the detections that were transformed back on to the original omnidirectional image. A different approach by Masato *et al.* [14] tries to train a rotation invariant model by introducing rotation augmentation during training, to partially overcome the rotation distortion of omnidirectional images. Both works face a similar challenge with only a limited amount of available omnidirectional data. To overcome this challenge, they either only work on the unwrapped image to better fit the model dataset. Or by rotating the large datasets, to better resemble how people appear on the omnidirectional images. Our application tries to preserve the privacy by using low-resolution image resolution. Unwrapping these low resolution images or using rotated low-resolution images will

not use the environment specific data to compensate the loss in data. By training models on the low-resolution omnidirectional data, we expect the model can better learn what describes the low-resolution representation of people.

Previous work by Callemeyn *et al.* [3] follows this methodology and faces a similar challenge. They propose an approach to count the number of people present in meeting-rooms and flex desk environments, while working towards privacy preservation. They also have a limited amount of omnidirectional data suited to their use-case, and therefore recorded their own data in several scenarios. Since this data was unlabelled, they proposed a teacher-student approach, where the teacher uses Yolov2 [2] and OpenPose [15], [16] detector models on unwrapped images to first generate labels on their private dataset.

Based on these generated labels they train several Yolov2 models and increase the privacy by reducing the image resolution. By lowering the resolution they reduce the details that make a person recognisable, increasing the privacy. However, their teacher pipeline was optimised for their omnidirectional camera and produced large area detections. Furthermore, they reached a architectural low-resolution limit and only went as low as  $96 \times 96$  pixels. Our approach proposes flexible detection generation pipeline that produces smaller annotations. We decreased the resolution further, and propose a novel approach that uses temporal data to retain performance.

### III. APPROACH

We propose a two part approach, capable of counting the amount of people on privacy preserving low resolution omnidirectional data. The omnidirectional camera will capture a static overview of meeting rooms or flex-desks. The only data variation is caused by the present people in the room. We therefore suggest to train a specific model on each scene, instead of generic for multiple scenes. However, recording new data for each scene requires manual data labelling before the data can be used for training. To significantly decrease the needed amount of manual labour, we propose an approach capable of autonomously annotating the data, described in

section III-A. After autonomously acquiring bounding box labels on the high resolution data, we train several models for extremely low resolutions using these labels. Lowering the resolution will decrease the image detail and increase the sense of privacy. However, decreasing image resolution also leads to a significant loss in spatial data. We therefore propose an approach that is able to retain the model performance using temporal information, described in section III-B.

#### A. Generating Bounding boxes

The high-resolution omnidirectional input image is first unwrapped into  $k$ -images with *overlap* at either side. Figure 1 shows an example with  $k = 3$  and an overlap of 10%. Additionally we also determine exclusion areas, for example near the heavily distorted centre or upper boundary where people will never be present. These parameters ( $k$ , *overlap*,  $y_b$ ) determines the number of unwrapped fragments and the width and height of each fragment. Each fragment will be used as input for both the Yolov2 [2] person detector and OpenPose [16] pose estimator. Since the unwarping parameters greatly influence the performance of the second step, the optimal settings were determined experimentally, as discussed in section IV-A.

To warp the detection out on the omnidirectional image, we first transform each bounding box to a poly-point representation. Instead of only using the bounding box corners, we add 2 evenly separated points on each of the bounding box sides. The upper left and upper right corners of the bounding box warped on the omnidirectional image will be placed further away from each-other. When a new bounding box is calculated based on these warped detection points, the top corners will enlarge the detection area greatly. By removing the top corners of each detection before warping, our warped detection will have a smaller area that better fits around the person.

Both the calculated Yolov2 points and the OpenPose pose estimation output are then warped back on the omnidirectional image. Around each set of point-detections we calculate a bounding box, and suppress overlapping detections using NMS with a threshold of 0.4.

#### B. Interlacing kernel

Previous work by Callemein *et al.* [3] shows it is possible to detect people in similar scenes even after decreasing the image resolution. However, they are only able to reduce the resolution to 96 pixels (due to architecture limits). At such resolutions, people arguably remain recognisable. Our approach allows for extreme low resolutions, not limited by network constraints.

We use the Yolov2 object detection architecture to train several models, using the autonomously generated bounding boxes, discussed in section III-A, with a network resolution of 160px and 96px. Yolov2 uses network resolution resize augmentation to allow the model to learn different scales. For this purpose, they resize their network between a range of [320; 608] with a step of 32. We follow a similar approach and allow the network to randomly resize the network resolution within a range of  $[net_{res} - 32*2; net_{res} + 32*2]$ . Architectural

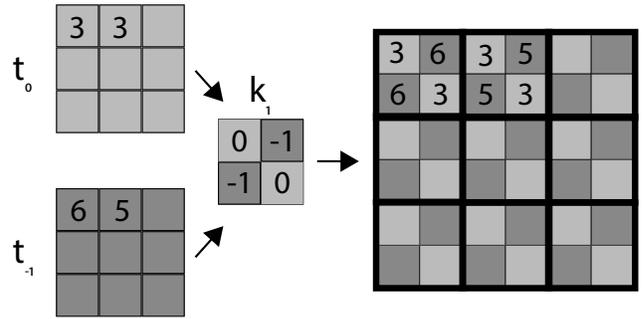


Fig. 4: overview showing an example temporal upscale combining multiple frames, resembling interlacing

limitations only allow the lowest network resolution of 96px. In the case of  $net_{res} = 96$ , the network will only have random upscales and saturate at 96px. While  $net_{res} = 160$ , is still able to use the full random resize scope.

Apart from decreasing the network resolution, we propose several down-up-scale resolutions (64px, 48px and 32px). The images are first down-scaled to these extreme low resolutions and up-scaled with linear interpolation to the network resolution. For each network resolution and down-scale resolution a model was trained to assess its performance.

While reducing the resolution further increases the privacy, it will also lead to general detail and data loss, increasing the challenge to detect people. To cope with this, we expect the region around people to contain temporal information, on the assumption that people constantly move over time, while e.g. furniture remains static. Based on this assumption, we propose to use this temporal information by combining multiple low-resolution images with *interlacing kernels*. Interlacing, used in video compression and computer graphics, is a technique where only parts of each frame are stored (reducing the required storage) and shown when watching the video. These interlaced frames combine multiple low-resolution frames together into a single higher resolution image. In our case we use a similar concept as interlacing, namely combining temporal-spatial information to increase the image resolution. This way we can improve performance compared to a linear interpolation up-scale. We therefore propose to use kernels (further referenced as *interlacing kernels*) to combine multiple frames together. Figure 4 illustrates our proposed approach, combining two  $3 \times 3$  matrices into a  $6 \times 6$  matrix using a  $2 \times 2$  interlacing kernel. This interlacing kernel contains index values and determines the source matrix to gather pixel data from. When no movement has taken place, these pixels will show similar behaviour as up-scale interpolation. In case of movement, we assume that the pixels will contain this movement and show more edges opposed to interpolation output.

## IV. EVALUATION

To evaluate both our autonomous label generation, described in section III-A, and our proposed approach to combine spatial-temporal data, described in section III-B, we use two

Dataset		Frames	Sequence	Dataset Annotation type	People	Level of movement
Mirror Challenge	train a	1084	2, 3	public bounding boxes	3	moderate
	train b	3608	7,8,10,12,13,14,17,18			
	test a	1123	1,4,5	manually annotated bounding boxes		
	test b	2246	9,11,15,16			
PIROPO	train a	10969	omni_1a	head points	3	high
	train b	4585	omni_1b			
internal office	train	7686		None	4	limited
	test	9953		manually annotated bounding boxes		

TABLE I: Used datasets during our experiments, showing the aggregated sequences, the number of frames, people and movement level.

publicly available datasets, PIROPO<sup>1</sup> and MirrorChallenge<sup>2</sup>. Both the PIROPO and MirrorChallenge dataset contain multiple camera setups and positions, we only use the open space and flex desk sequences since they better fit our case. In order to further test our system and simulate real office space situations, we recorded a private office dataset with little movements since the people are at their desks. Table I shows the summary of used datasets during our experiments.

Section IV-A evaluates the proposed approach to generate bounding box labels on all three training datasets. Based on the best settings, we will then use these automatically generated labels to train several models, evaluated on the test datasets in section IV-B.

#### A. Automatic labelling

We can only evaluate our automatic bounding box generation technique on both the MirrorChallenge and PIROPO training datasets, since our private office dataset has no training labels. The PIROPO dataset, however, has no bounding box annotations, only head point annotations. We therefore use point-wise evaluation and checking whether the detection box contains the head annotation point. When this is the case, we annotate it as a true positive, when the point is outside any of the detection boxes it is counted as a false negative. The remainder of the detections that was not matched with a head annotation are counted as false positives.

Figures 5 and 6 illustrate two precision-recall curves for the training set A and B of the MirrorChallenge dataset. The leftmost pr-curve shows the results when using point-wise evaluation, the rightmost will compare the bounding boxes with an IoU of 0.4. As mentioned in section III-A different settings can be used to generate the bounding boxes. In our case we evaluated different amount of k-frames, with  $k = [2; 3; 4]$  for YoloV2 and  $k = [2; 3]$  for OpenPose. Since the PIROPO training datasets only have point annotations, figure 7 only illustrates the pr-curves using point evaluation, with  $k = [2; 3; 4; 5]$  for YoloV2 and  $k = [2; 3]$  for OpenPose. On the PIROPO train set A, we noticed that most of the head annotations were near the circular boundary of the omnidirectional image. We therefore set  $y_b$  to only use the upper-part when unwarping the omnidirectional image,

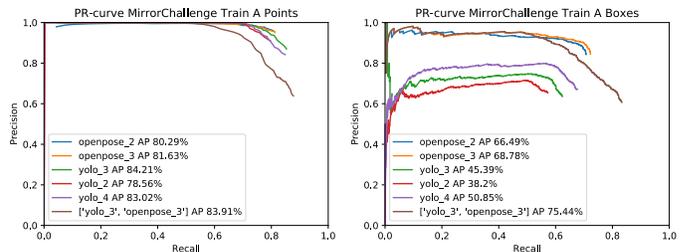


Fig. 5: Mirror Train A generated annotations PR-curves evaluating the boxes with an IoU=0.4 and point within bounding box.

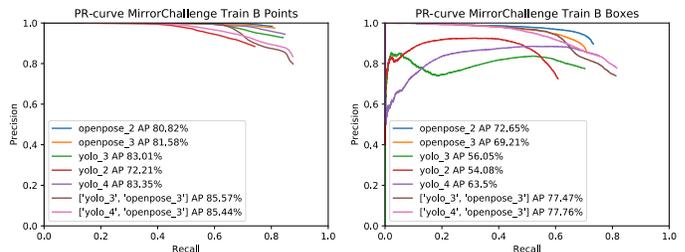


Fig. 6: Mirror Train B generated annotations PR-curves evaluating the boxes with an IoU=0.4 and point within bounding box.

producing fragments with a large width and a small height. By increasing  $k$ , we improve the width/height ratio to better fit our detection architectures, resulting in higher accuracy. The best performing OpenPose and YoloV2 detections are then combined using NMS. To further increase the accuracy, we compare the number of detections with the mean number of detections of the past. If the current frame has a number of detections not equal to the mean number, we allow the system to drop the current frame. We then determine the optimal threshold with the F1 score and use this threshold to generate annotations that were used for training the models described in section III-B.

To test whether the automatically generated annotations are adequate enough, we trained three models with different network resolutions [448; 160; 96] for each training dataset that were evaluated on the test sets. In table II you can find the average precisions, showing good results on the original 448 resolution, and a slight decrease after decreasing the network resolution. This shows that we are capable of training models

<sup>1</sup><https://www.gti.ssr.upm.es/research/gti-data/databases>

<sup>2</sup><https://www.hcd.icat.vt.edu/mirrorworlds/challenge/index.html>

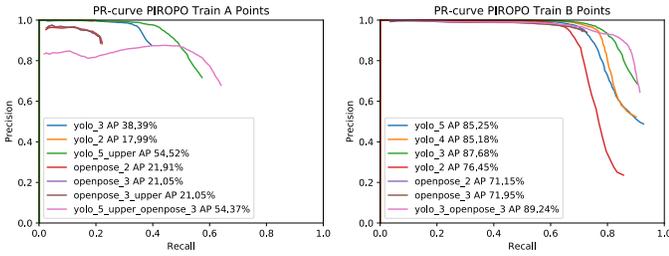


Fig. 7: PIROPO Train A and B generated annotations PR-curves evaluating point within bounding box.

Dataset	Net resolution		
	448	160	96
PIROPO Test B	0.966	0.817	0.428
MirrorChallenge Test A	0.878	0.670	0.676
MirrorChallenge Test B	0.942	0.911	0.864
Private Office	0.967	0.831	0.791

TABLE II: Average precision results for each network input resolution, trained on the training sets with automatically generated labels and evaluated on the test sets.

with our generated annotations with acceptable performances.

### B. Interlacing kernel

The main focus of this paper is counting the people present in omnidirectional images, while lowering the resolution to increase the privacy preservation. Note that our purpose is to make people detectable, but not identifiable. Thus the absolute size of the pixel has less importance than the relative pixel size to the size of the face. However, in our case, we use fixed ceiling mounted top down looking cameras with a wide angle lens adding lens distortion. This implies that the size

Dataset	Net Res	Scale res	Linear	Interlacing Kernel		
				k1	k2	k3
PIROPO TEST B	160	64	0.746	0.824	<b>0.860</b>	0.655
		48	0.745	0.757	<b>0.833</b>	0.750
		32	0.759	0.774	<b>0.853</b>	0.748
	96	48	0.351	0.428	<b>0.526</b>	0.323
		32	0.271	<b>0.306</b>	0.232	0.135
MIRROR TEST A	160	64	0.575	<b>0.758</b>	0.648	0.717
		48	0.636	<b>0.806</b>	0.686	0.698
		32	0.585	<b>0.673</b>	0.552	0.464
	96	48	0.435	0.191	<b>0.610</b>	0.126
		32	0.557	0.169	<b>0.707</b>	0.257
MIRROR TEST B	160	64	<b>0.939</b>	0.897	0.885	0.866
		48	<b>0.926</b>	0.895	0.906	0.864
		32	0.916	0.882	<b>0.917</b>	0.877
	96	48	<b>0.877</b>	0.859	0.869	0.814
		32	<b>0.872</b>	0.850	0.838	0.852
PRIVATE TEST	160	64	<b>0.866</b>	0.651	0.793	0.797
		48	0.751	0.825	<b>0.879</b>	0.641
		32	<b>0.710</b>	0.520	0.616	0.439
	96	48	0.641	0.583	0.639	<b>0.669</b>
		32	<b>0.66</b>	0.583	0.652	0.392

TABLE III: PIROPO, MIRROR A, MIRROR B and Office average precisions on different network and down-scale resolution comparing linear upscale vs. upscale with interlacing kernels.

Net Res	Scale Res	Linear	Interlacing Kernel with delta t		
			1	2	3
160	64	0.866	0.793	<b>0.878</b>	0.844
	48	0.751	0.879	0.794	<b>0.892</b>
	32	0.710	0.616	<b>0.7836</b>	0.741

TABLE IV: Results on the private dataset, using kernel  $k_2^t$  for  $t = [1; 3]$ .

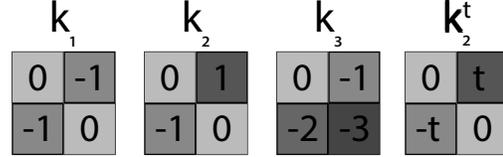


Fig. 8: Proposed interlacing kernels used to up-scale the low-resolution images.

of the persons relative to the overall camera image is already small. During evaluation we will use point evaluation better evaluating the counting output, while this still takes in account the rough location of the detection.

Table III show the average precisions for all trained models on each dataset, with a network resolution of 160 and 96, on images that were down-scaled to a smaller resolution. We compare linear interpolation with three temporal interlacing up-scaling kernels  $k_1; k_2; k_3$ . The three leftmost kernels illustrated in figure 8 were used, where the value represents the frame index time difference to be used to upscale the images, the datasets are all recorded at 15 fps. A time-delta of 1 will result in a 66ms shift in time. We observe that the effect of taking time into account by using these interlacing kernels depends on the amount of movement of the people in the images. People walking around indeed will generate much more interlacing artefacts than people sitting immobile at their desks. The results on the PIROPO test sets, with high level of movement, Table III shows that in all cases this model outperforms up-scaling the low-resolution images using linear interpolation. On the MirrorChallenge test sets, with a moderate level of moment, the interlaced upscale does not always outperform the linear interpolation models, but shows similar results. On our private office dataset, with a very limited level of movement, the same conclusion as on the MirrorChallenge counts, and both perform well. Yet, the interlaced up-scale in some cases shows a little increase or decrease in performance.

On our private office dataset a smaller impact of time with kernel  $k_2$  is visible and to be expected, since little variation occurs in a 132ms time window. The kernel base-structure as  $k_2$  was used, adjusted to hold a variable time-delta  $t$ , kernel  $k_2^t$ , illustrated in the rightmost kernel in figure 8. The results for models trained with  $k_2^t$  with  $t = [1, 2, 3]$  are show in table IV, showing that increasing the time-delta increases the performance and outperforms the linear interpolation which was used as the baseline.

On our private dataset, table II shows the average precisions

Device	Resolution	Seconds per Frame
Raspberry Pi 2	448	18.60
	160	3.60
	96	2.17
Raspberry Pi 3B	448	16.60
	160	2.96
	96	1.83
Raspberry Pi 3B+	448	11.72
	160	2.07
	96	1.30

TABLE V: Processing speed of our models on embedded platforms.

using our baseline, which is the approach of Callemain *et al.* [3]. Table IV shows the results with our proposed approach on the same dataset. We indeed see that our newly proposed approach enables to reduce the image resolution drastically further (three times lower as in [3]), while keeping the average detection precision similar. Figure 9 shows example images from the each test set, showing the ground truth (red) and the low-resolution based detections (blue) on both original high-resolution image (left) and the low-resolution image (32px with interlacing kernel  $k_2$ ).

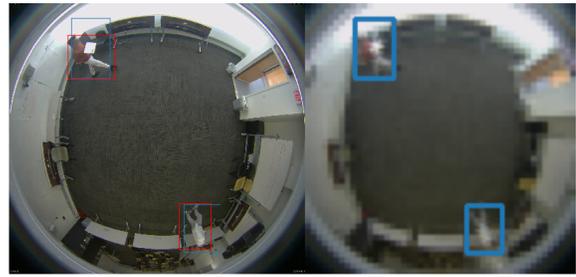
### C. Embedded implementation

We implemented the resulting networks on Raspberry Pi 2, 3, and 3B devices, sporting a RaspiCam camera with a 1.1 mm omnidirectional lens. Table V shows the measured frame rates achieved by our models on these embedded platforms, after automatic self-training.

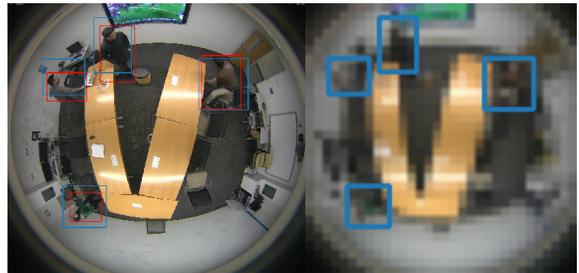
A measurement update rate of 1.3s to 4s does maybe not seem enormous, but undeniably it is very suited for the application at hand. The number of people in a room must certainly not be measured more frequently for such a room reservation system.

## V. CONCLUSION

In this paper, we presented an omnidirectional camera-based sensor counting the number of people in flex-desk and meeting room environments. To overcome the scarce amount of labeled omnidirectional data, an autonomous label generation system based on state-of-the-art person detectors was proposed, allowing for scene-specific data recording, label generation and training of several models. By decreasing the image resolution, we achieve true privacy-preservation, reducing the input image resolution to the utmost. To retain high detection accuracy, we proposed to incorporate temporal data to compensate for the loss in spatial data. Results showed that our approach is capable of using scene knowledge to generate labels that can be used during training. Evaluating the models, trained on these generated labels, showed that our generated labels were adequate enough to train models capable of counting people with high accuracy. Furthermore, by using interlacing kernels that take into account temporal information, we see a clear improvement over normal interpolation up-scaling techniques.



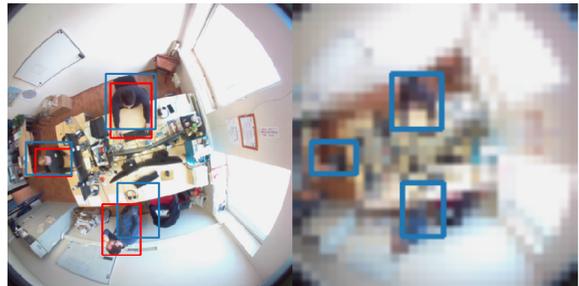
(a) MIRROR TEST A



(b) MIRROR TEST B



(c) PIROPO TEST B



(d) PRIVATE TEST

Fig. 9: Dataset examples, showing the annotations (red) and detections (blue) on both the high and low resolution frames (Net Res: 160 Scale Res: 32 with interlacing kernel  $k_2$ ).

## REFERENCES

- [1] D. J. Butler, J. Huang, F. Roesner, and M. Cakmak, "The privacy-utility tradeoff for remotely teleoperated robots," in *ACM/IEEE International Conference on Human-Robot Interaction*, pp. 27–34, ACM, 2015.
- [2] J. Redmon and A. Farhadi, "Yolo9000: Better, faster, stronger," *arXiv preprint arXiv:1612.08242*, 2016.
- [3] T. Callemain, K. Van Beeck, and T. Goedemé, "How low can you go? privacy-preserving people detection with an omni-directional camera," in *International Joint Conference on Computer Vision, Imaging and*

- Computer Graphics Theory and Applications.*, International Joint Conference on Computer Vision, Imaging and Computer , 2018.
- [4] C. Arteta, V. Lempitsky, and A. Zisserman, "Counting in the wild," in *European conference on computer vision*, pp. 483–498, Springer, 2016.
- [5] Y. Zhang, D. Zhou, S. Chen, S. Gao, and Y. Ma, "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, 2016.
- [6] D. B. Sam, S. Surya, and R. V. Babu, "Switching convolutional neural network for crowd counting," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4031–4039, IEEE, 2017.
- [7] J. Zou, Q. Zhao, W. Yang, and F. Wang, "Occupancy detection in the office by analyzing surveillance videos and its application to building energy conservation," *Energy and Buildings*, vol. 152, pp. 385–398, 2017.
- [8] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, pp. 91–99, 2015.
- [9] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*, pp. 21–37, Springer, 2016.
- [10] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- [11] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results." <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- [12] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*, pp. 740–755, Springer, 2014.
- [13] R. Seidel, A. Apitzsch, and G. Hirtz, "Improved person detection on omnidirectional images with non-maxima suppression," *arXiv preprint arXiv:1805.08503*, 2018.
- [14] M. Tamura, S. Horiguchi, and T. Murakami, "Omnidirectional pedestrian detection by rotation invariant training," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1989–1998, IEEE, 2019.
- [15] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *CVPR*, 2017.
- [16] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, "OpenPose: realtime multi-person 2D pose estimation using Part Affinity Fields," in *arXiv preprint arXiv:1812.08008*, 2018.