

Privacy-Preserving Person Detection Using Low-Resolution Infrared Cameras

Thomas Dubail, Fidel Alejandro Guerrero Peña, Heitor Rapela Medeiros,
Masih Aminbeidokhti, Eric Granger, and Marco Pedersoli

Laboratoire d'imagerie, de vision et d'intelligence artificielle (LIVIA)
Dept. of Systems Engineering, ETS Montreal, Canada
{thomas.dubail.1, fidel-alejandro.guerrero-pena.1,
heitor.rapela-medeiros.1, masih.aminbeidokhti.1}@ens.etsmtl.ca,
{eric.granger, marco.pedersoli}@etsmtl.ca

Abstract. In intelligent building management, knowing the number of people and their location in a room are important for better control of its illumination, ventilation, and heating with reduced costs and improved comfort. This is typically achieved by detecting people using compact embedded devices that are installed on the room's ceiling, and that integrate low-resolution infrared camera, which conceals each person's identity. However, for accurate detection, state-of-the-art deep learning models still require supervised training using a large annotated dataset of images. In this paper, we investigate cost-effective methods that are suitable for person detection based on low-resolution infrared images. Results indicate that for such images, we can reduce the amount of supervision and computation, while still achieving a high level of detection accuracy. Going from single-shot detectors that require bounding box annotations of each person in an image, to auto-encoders that only rely on unlabelled images that do not contain people, allows for considerable savings in terms of annotation costs, and for models with lower computational costs. We validate these experimental findings on two challenging top-view datasets with low-resolution infrared images.

Keywords: Deep Learning, Privacy-Preserving Person Detection, Low-Resolution Infrared Images, Weak Supervision, Embedded Systems.

1 Introduction

Intelligent building management solutions seek to maximize the comfort of occupants, while minimizing energy consumption. These types of solutions are crucial for reducing the use of fossil fuels with a direct impact on the environment. Such energy-saving is usually performed by adaptively controlling lighting, heating, ventilation, and air-conditioning (HVAC) systems based on building occupancy, and in particular the number of people present in a given room. For this, low-cost methods are needed to assess the level of room occupancy, and efficiently control the different systems within the building.

Among the different levels of occupancy information that can be extracted in an intelligent building [25], Sun et al. define the location of occupants as the most important and fine-grained for smart building control. Given the recent advances in machine learning and computer vision, most solutions usually rely on deep convolutional neural networks (CNNs) to detect people [8,7]. Despite the high level of accuracy that can be achieved with CNNs for visual object detection based on RGB images, their implementation for real-world video surveillance applications incurs in a high computational complexity, privacy issues, and gender and race biases [4,22]. Finally, building occupancy management solutions are typically implemented on compact embedded devices, rigidly installed on the ceiling or portals of rooms, and integrating inexpensive cameras that can capture low-resolution IR images.

To mitigate these issues, He et al. [9] have proposed a privacy-preserving object detector that blurs people’s faces before performing detection. To strengthen the detector against gender/race biases, the same authors proposed a face-swapping variation that also preserves privacy at the cost of increased computational complexity. Regardless of the good performance, their approach does not ensure confidentiality at the acquisition level, relying on RGB sensors to build the solution. Furthermore, their detector was designed for fully annotated settings using COCO [16] as the base dataset. This makes it difficult to generalize to people detection under different capture conditions (like when cameras are located on the ceiling on compact embedded systems), and extreme changes in the environment. In addition, it is difficult to collect and annotate image data to train or fine-tune CNN-based object detectors for a given application, so weakly-supervised or unsupervised training is a promising approach.

In contrast, our work tackles the occupants location problem by detecting people in infrared (IR) images at low resolution, which avoids most of the above-mentioned issues on privacy. Low resolution not only reduces computational complexity but also improves privacy, i.e., a detection on high-resolution infrared images would not be enough as it is possible to re-identify people [32]. More specifically, we analyse people detection with different levels of supervision. In this work, we compare unsupervised, weakly-supervised and fully supervised solutions. This is an essential aspect of the detection pipeline since producing bounding box annotations is very expensive, and there is a lack of good open-source object detection datasets for low-resolution infrared scenarios. In fact, reducing the level of supervision can lead to improved scalability for real applications and reduced computation, which is important considering the use of the proposed algorithms on embedded devices.

The contributions of this work are the following. (i) We propose cost-effective methods for estimating room occupancy under a low-supervision regime based on low-resolution IR images, while preserving users privacy. (ii) We provide an extensive empirical comparison of several cost-effective methods that are suitable for person detection using low-resolution IR cameras. Results indicate that, using top-view low-resolution images, methods that rely on weakly-labeled image data can provide good detection results, and thus save annotation efforts

and reduce the required complexity of the detection model. (iii) To investigate the performance of person detecting methods on low-resolution IR images, our results are shown in two challenging datasets – the FIR-Image-Action and Distech-IR datasets. Finally, we provide bounding box annotations for the FIR-Image-Action [31] dataset¹.

2 Related Work

Privacy-preserving methods are of great interest to the scientific community. As a result, multiple approaches have been proposed in the past to circumvent the challenge. Ryoo et al. [21] proposed a method for learning a transformation that obtained multiple low-resolution images from a high-resolution RGB source. The method proved effective for action recognition even when inputs’ resolution were down-sampled to 16×12 . These findings were validated by others [29,6] using similar downsampling-based techniques to anonymize the people displayed in the images. On the other hand, recent approaches [20,9] have focused on producing blurred or artificial versions of people’s faces while preserving the rest of the image intact. These methods usually rely on Generative Adversarial Networks (GANs) to preserve image utility, while producing unidentifiable faces. Specifically, the method of He et al. [9] is one of the first approaches to apply such anonymization in an object detection task. Despite recent advances in the area, these proposals are specialized for RGB images, which are anonymized after the undisclosed acquisition. As an alternative to RGB cameras, others authors [26,15,27,28] have proposed using low-resolution IR cameras to preserve anonymity at the phase of image acquisition. While most of these works target action recognition task, we focus on people detection as in [24,5].

Complementary to privacy preservation, this work also focuses on studying detection methods with different levels of supervision. Such a study aims to find techniques that reduce annotation costs without a significant performance reduction when compared with fully supervised approaches. In this work, we followed the auto-encoders (AE)-based anomaly-detection method similar to the one proposed by Baur et al. [2]. The technique is used to learn the distribution of typical cases, and applied later to identify abnormal regions within the image. However, different from their proposal, we focus on object detection instead of image segmentation. We also evaluate weakly-supervised detection methods based on Class Activation Maps (CAMs) [33,23] following the authors algorithm. Nonetheless, we use a customized CNN to be consistent with our low-resolution inputs. Finally, we consider also the fully supervised techniques Single Shot Detector [17] and Yolo v5 [11] as upper-bound references, sharing the same backbone as the previously described approaches. The aim of this work is not to use the latest developments for each type of technique. Instead, we aim to achieve a good trade-off between performance and computational complexity with simple and commonly used techniques. In particular, we favor simple occupants’ de-

¹ <https://github.com/ThomasDubail/FIR-Image-Action-Localisation-Dataset>

tection approaches that can run in embedded devices and with capabilities for handling low-resolution infrared images.

3 Person Detection with Different Levels of Supervision

Several cost-effective methods may be suitable for person detection using low-resolution IR cameras installed on ceilings. The goal of object detection is to find a mapping f_θ such that $f_\theta(x) = z$, where z are the probabilities that a bounding box belongs to each class. Note that such a mapping can be obtained using any level of supervision. In this paper, we seek to compare the detection accuracy of methods that rely on different levels of image annotation, and thereby assess the complexity needed to design embedded person detection systems.

In this work we consider a “fully annotated” dataset of IR images at low resolution $\mathcal{F} = \{(x_0, b_0), (x_1, b_1), \dots, (x_N, b_N)\}$, with x an IR image and $b = \{(c_0, d_0, w_0, h_0), \dots, (c_B, d_B, w_B, h_B)\}$ rectangular regions enclosing the objects of interest, also known as bounding boxes. Without loss of generality, we use a center pixel representation (*center x, center y, width, height*) [17] for defining the bounding boxes. In the given formulation, all bounding boxes belong to the persons’ category. Consequently, a “weakly annotated” IR dataset is defined as $\mathcal{W} = \{(x_0, y_0), \dots, (x_N, y_N)\}$ in which $y_i \in \{0, 1\}$ corresponds to an image-level annotation indicating whether a person is present ($y_i = 1$) or not ($y_i = 0$) in the image x_i . Finally, at the lowest level of supervision, an “unlabeled” dataset containing only IR images without annotations is expressed as $\mathcal{U} = \{x_i\}$. Please note that for this study, the datasets \mathcal{F} , \mathcal{W} , and \mathcal{U} are drawn from the same pool of IR images but with different levels of annotations.

The rest of this section details the methods compared in this paper, each one trained according to a different level of supervision. Here, the backbone of all deep learning based methods remained the same. We focus on low-cost methods that can potentially be implemented on compact embedded devices.

3.1 Detection through thresholding

Let x be an IR thermal image from \mathcal{U} . The people within the images appear as high-temperature blobs easily distinguishable from the low-temperature background, Figure 1a. Such a property allows us to directly apply a threshold-based mapping $g_\tau(x) = \Phi(\llbracket x \geq \tau \rrbracket)$ to obtain the persons’ location. In this formulation, we use $\llbracket \cdot \rrbracket$ to refer to the Iverson bracket notation, which denotes the binarization of the image x according to the threshold τ , Figure 1b. Here, the value of τ is manually determined according to a validation set, or automatically following Otsu’s method [19], hereafter referred as *Threshold* and *Otsu’s Threshold* respectively. Finally, a mapping Φ converts the segmentation map into a bounding box by taking the minimum and maximum pixels from each binary blob (see Figure 1c). This method is more seen as a post-processing step for the following methods, although we have evaluated it to have a lower bound for detection.

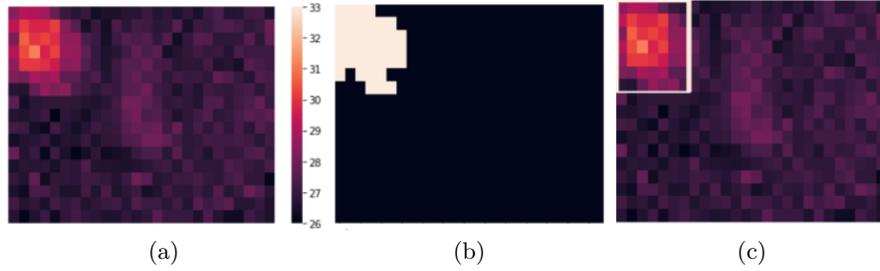


Fig. 1: Example of an IR image (a) that is binarized using the threshold $\tau = 29$ (b), and represented as a bounding box (c).

3.2 Unsupervised anomaly-based detection using auto-encoders

For this next approach, people are considered anomalies for the distribution of empty rooms. In this work, we follow the method proposed by Baur et al. [2] to model the background distribution using an auto-encoder (AE) f_θ trained using only empty rooms, $\mathcal{W}^0 = \{x_i \mid y_i = 0\}$. Such an approach acts as a background reconstruction technique whenever an anomaly is present, i.e., the AE will not be able to reconstruct it. Then, we can highlight the anomaly by taking the difference between the input image and the obtained reconstruction, $x - f_\theta(x)$ (see Figure 2). Finally, the anomaly detection method for person detection is defined as $\Lambda_{\theta, \tau}(x) = g_\tau(x - f_\theta(x))$ where g_τ is the thresholding technique explained in the previous section. Thus, the detection is performed in a two-step process: anomaly boosting and anomaly segmentation-localization, which can be done by setting a threshold τ .

The encoder architecture comprises six convolutional layers with kernels of size 3×3 . Max-pooling operations are used every two convolutional layers to increase the field of view. The decoder follows a symmetrical architecture of the encoder using transposed convolution with a stride of 2 as upsampling technique. The bottleneck uses a linear layer with 256 neurons which encode input infor-

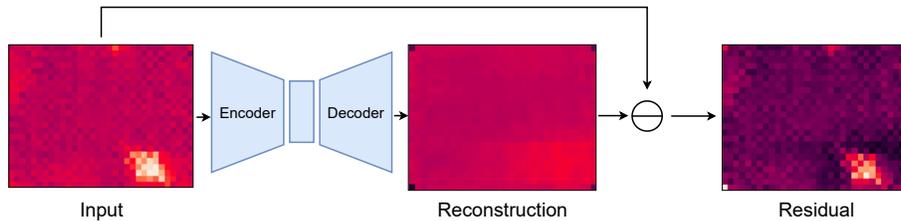


Fig. 2: Unsupervised anomaly-based method for person detection from low-resolution IR images.

mation as a vector projected onto the latent space. Finally, a reconstruction loss is used to guide the training process and find a feasible minimizer θ for solving our background reconstruction task. In this work the Mean Square Error (MSE) loss is used, $\mathcal{L}_{MSE}(x, \theta) = \frac{1}{|W^0|} \sum_{x_i \in W^0} (x_i - f_\theta(x_i))^2$. In later sections we refer to this approach as *deep auto-encoder* or simply *dAE*.

Besides the classical AE, several types of hourglass architectures have been used for anomaly detection [2]. One of the most popular versions are the *variational auto-encoders (dVAE)* [13] where the latent vector is considered to be drawn from a given probabilistic distribution. Here we use the same architecture for both AE-based methods with a distinction for the loss function where the KL-divergence regularization is added to the MSE loss to enforce a normal distribution to the latent space.

3.3 Weakly-supervised detection using class activation mapping

Let (x, y) be a generic tuple from \mathcal{W} where x is an image and $y \in \{0, 1\}$ is a category indicating whether a person is present or not in x . A weakly-supervised approach for object localization is such that, by exploiting only the image-level annotation during training, learns a mapping c_φ to retrieve the object location during the evaluation. This kind of approaches based on Class Activation Maps (CAM) techniques [30,14,1] has been widely explored in the literature. The principle is based on the use of the compound function $c_\varphi(x) = (c_\psi^1 \circ c_\phi^0)(x)$ that relies on a feature extractor c_ϕ^0 followed by a binary classifier c_ψ^1 . Here c_ϕ^0 is implemented as a CNN and c_ψ^1 as a Multi-layer Perceptron. Then, the following minimization based on cross-entropy loss is performed in order to find the optimal set of weights: $\min_\theta - \sum_{(x,y) \in \mathcal{W}} y \cdot \log c_\varphi(x)$. Once a feasible set of weights is found, a non-parametric transformation function uses the output from the feature extractor c_ϕ^0 to produce an activation map for each category in the task, $M(c_\phi^0(x))$. Note that in this task, the computation of such an activation map is only performed whenever a positive classification is obtained, i.e., $c_\varphi(x) \geq 0.5$. The architecture used for c_ϕ^0 is the same as in the encoder for the *dAE* technique, but without using the Max-pooling layers to avoid losing resolution. In this work, we use three variants for the Global Average Pooling M . The first is the classic weighting approach proposed by [33], hereafter referred to as *CAM*, the second is the gradient-based CAM proposed by [23], known as *GradCAM*, and the last one is the hierarchical approach known as *LayerCAM* [10]. As in the previous techniques, the final localization is obtained using the thresholding-based mapping g_τ .

3.4 Fully-supervised detection using single shot detectors

At the higher level of supervision, we explore the mapping function with bounding box annotations within the cost function. Let h_ϑ be a mapping parameterized over ϑ , which produces bounding box predictions. Among the different types of

detectors existing in the literature, we used Single Shot Detector (*SSD*) [17] and *Yolo v5* [11] solutions since they are suitable for low-resolution images and allows using a custom backbone without serious implications for the training process. In this study, we enforce the same architecture for feature extraction as in the AEs and CAMs approaches. However, unlike the previous techniques, such supervised mappings do not require the composition with the thresholding function g_τ . Let $(x, b) \in \mathcal{F}$ be an IR image with its corresponding bounding box annotation. The learning process for both methods solves the optimization problem $\vartheta^* = \arg \min \mathcal{L}(x, b, \vartheta)$, by minimizing the cost of the model \mathcal{L} . Despite their differences in terms of representation, in both cases the loss function uses a supervised approach that measures the difference between the output $z = h_\vartheta(x)$ and the expected detections b .

4 Experimental Methodology

4.1 Datasets

In this study, two datasets were used to assess the IR person detection using models trained with different levels of annotation – the public FIR-Image-Action [31] dataset, and our Distech-IR dataset.

1) FIR-Image-Action with bounding box annotations

The FIR-Image-Action [31] dataset includes 110 annotated videos. We randomly selected 36 videos from this pool for the test and the others 74 for training and validation. Furthermore, training and validation sets were separated using a random selection of the frames (70% and 30%, respectively). All the approaches have been trained using the same data partition to ensure comparability.

To the best of our knowledge, there are no low-resolution IR datasets with bounding box annotations for person detection. Therefore, we annotated this dataset at bounding box level. The dataset was created by Haoyu Zhang of Visiongo Inc. for video-based action recognition. Such a dataset offers 126 videos with a total combined duration of approximately 7 hours. Since this study aims to evaluate the performance of different techniques for IR-based people localization, we only used the IR images provided by the authors for our experiments. Nevertheless, it is worth mentioning that two modalities are available within the dataset: RGB with a spatial resolution of 320×240 acquired at 24 FPS, and IR with a resolution of 32×24 and sampled at 8 FPS. Although the RGB falls outside the scope of this work, we used them for obtaining bounding box annotations, as described later. As part of this work’s contributions, we have publicly created and released the localization annotations for 110 videos out of the 126 for both RGB and IR modalities. Since there is redundancy within neighboring frames and our application does not require video processing, we further sampled the IR dataset obtaining the equivalent of 2 FPS videos.

We used a semi-automated approach to obtain bounding box annotations for the challenging low-resolution IR images in FIR-Image-Action. First, we create bounding box annotations by hand of a randomly selected subset of the RGB

frames. We carefully curated these bounding boxes to reduce the impact of the misplacement when decreasing the resolution for the IR modality. Then, an *SSD* detector [17], h_{ϑ} , was trained over the RGB annotated dataset, being used afterward to obtain pseudo-labels over the remaining unannotated partition of the RGB dataset. A new randomly selected subset is then curated and h_{ϑ} training is repeated but using a larger partition of the data. This process was repeated three times resulting in a fully annotated version of the RGB dataset.

Finally, bounding box annotations for the IR dataset are obtained by pairing images from both modalities, followed by a coordinate aligning procedure. Since the videos for IR and RGB were out of synchronization, the initial time shift was manually determined using an overlay visualization of both modalities (see Figure 3c). Such a synchronization was performed individually for every video. The final IR localization annotations were obtained by doing a linear interpolation of the bounding box coordinates from the labeled RGB dataset. The parameters for the alignment were estimated using linear regression. An example of the obtained bounding box annotation for both RGB and IR modalities can be observed in Figure 3a and b, respectively.

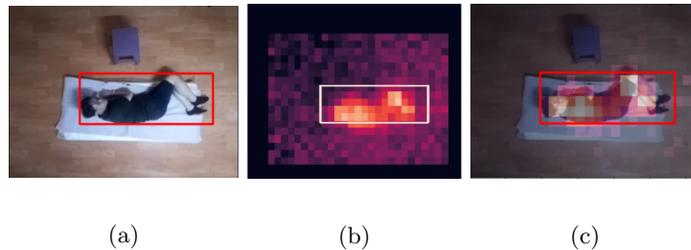


Fig. 3: Example of an RGB image with its ground truth (a), the corresponding IR image with the aligned bounding box (b), and an overlay of RGB and IR modalities (c) from the FIR-Image-Action dataset.

2) Distech-IR

The second dataset, named hereafter Distech-IR, followed the same separation proportions containing 1500 images for training, 500 for validation, and 800 for testing. Such a dataset, similar to FIR-Image-Action, contains two modalities of images (RGB and IR) with their corresponding bounding box level annotations provided by Distech Controls Inc. The dataset reflects the increasing interest by the industry for privacy preserving-based solutions for person localization and constitute an actual use case for this task. The Distech-IR dataset also proved to reflect better real-world scenarios since it is composed of seven rooms with different levels of difficulties, i.e., heat radiating appliances, sun-facing windows, and more than one person per room. For simulating deployment, we used rooms not seen before during training for the test set. Figure 4 shows some examples of images from both datasets.

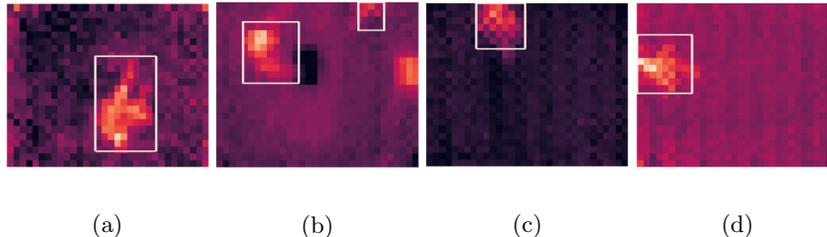


Fig. 4: Examples of IR images with their corresponding ground truth for FIR-Image-Action (a) and Distech-IR (b)-(d) datasets.

4.2 Implementation details

We used normalization by 50° to ensure small scale within the input map. Adam optimizer [12] was employed with an initial learning rate of 10^{-4} and decay of 0.2 with a patience of ten epochs. Additionally, a 15 epoch patience early-stopping was implemented. Then, the best model according to the validation loss was selected for each case. The time calculations were evaluated on an Intel Xeon CPU at 2.3 Ghz, however the training of the models was done on an NVIDIA Tesla P100 GPU. Each experiment was performed 3 times with different seeds. The validation protocols are presented with each dataset in section 4.1.

4.3 Performance metrics

In this study, we use the optimal Localization Recall Precision (oLRP) [18] in order to characterize the ability of each method to detect the presence of people and locate them. This metric allows us to evaluate with the same measurement methods that provide bounding box detections without a score associated (such as AEs and CAMs) and methods with a detection score (as SSD and Yolo v5 detectors). Additionally, as stated by the authors, it reflects the localization quality more accurately than other measurements, providing separate measures for the different errors that a detection method can commit. The metric takes values between 0 and 100, with lower values being better. As part of the metric computation, we also calculate the localization (oLRP_{loc}), False Positive (oLRP_{FP}), and False Negative (oLRP_{FN}) components which provide more insights of methods behavior. Finally, execution time on the same hardware has been computed to obtain an approximation of the time complexity.

5 Results and Discussion

Tables 1 and 2 summarize the obtained results for FIR-Image-Action and Distech-IR datasets, respectively. As expected, the fully supervised approaches obtained the best performance for most metrics and comparable results to the other approaches regarding False Negatives. The methods proved useful for locating people in diverse scenarios, even under low-resolution settings. However, they take

Table 1: Performance of detection methods on the FIR-Image-Action dataset. All metrics are calculated with an IoU of 0.5.

Model	oLRP ↓	oLRP_{loc} ↓	oLRP_{FP} ↓	oLRP_{FN} ↓	Time(ms) ↓
Threshold	86.5 ± 0.1	32.3	45.3	44.2	0.4
Otsu’s Threshold	83.5 ± 0.1	31.6	45.7	27.7	0.7
dVAE	74.7 ± 1.3	31.2	26.6	24.5	13.0
dAE	77.4 ± 1.3	30.4	31.2	29.1	12.3
CAM	85.1 ± 1.1	34.3	41.2	29.0	11.4
GradCAM	85.5 ± 3.2	34.5	43.0	32.1	24.3
LayerCAM	84.8 ± 2.2	34.9	37.2	33.1	25.6
SSD	63.8 ± 2.7	25.3	12.6	18.6	46.6
Yolo v5	56.9 ± 1.8	25.5	6.3	6.2	45.9

longer to execute than the second-best performed techniques, which is an important downside to take into account for measuring real-time occupancy levels in intelligent buildings. In particular, the *dAE* approach was 3.7 times faster than *Yolo v5* and 3.8 times faster than *SSD*.

We can refer to the AE-based anomaly detection approaches as second-placed strategies. Both *dAE* and *dVAE* showed similar performance in terms of LRP and efficiency between them. Furthermore, the techniques obtained localization performance comparable to the fully-supervised approaches, especially in more complex situations like the Distech-IR dataset. This result is remarkable considering that localization supervision is not used during AE training, and only empty room images were used (10% of the data in the FIR-Image-Action and 30% in Distech-IR). The methods also showed acceptable execution times for the application.

CAM methods provide a lower level of performance than other methods, despite having access to class-label annotations for training. Indeed, these methods

Table 2: Performance of detection methods on the Distech-IR dataset. All metrics are calculated with an IoU of 0.5.

Model	oLRP ↓	oLRP_{loc} ↓	oLRP_{FP} ↓	oLRP_{FN} ↓	Time(ms) ↓
Threshold	93.6 ± 2.4	37.1	72.7	54.1	0.4
Otsu’s Threshold	95.5 ± 1.2	34.2	83.3	50.0	0.7
dVAE	83.3 ± 8.9	33.2	32.7	40.6	13.0
dAE	82.7 ± 9.0	32.4	33.7	40.0	12.3
CAM	93.1 ± 1.9	37.6	59.7	52.3	11.4
GradCAM	91.6 ± 2.3	37.5	50.3	48.8	24.3
LayerCAM	91.1 ± 2.5	37.7	45.5	50.3	25.6
SSD	82.0 ± 7.2	31.1	26.3	44.7	46.6
Yolo v5	80.2 ± 7.7	30.2	31.4	37.4	45.9

are known to activate strongly for discriminant regions of an input image (since the backbone CNN is trained to discriminate classes), and be affected by complex image backgrounds [3]. These two factors affect its ability to define precise contours around a person.

As can be seen in the tables, *Otsu's Threshold* provides good person localization. However, it assumes a multi-modal intensity distribution for finding the threshold, which leads to false person localization in empty rooms. This effect can be observed by the high values of oLRP_{FP} . The rest of the approaches showed comparable results to this last one but were still far from the fully-supervised process. Figure 5 shows some examples of the obtained result over the FIR-Image-Action for *dVAE*, *gradCAM*, and *Yolo v5* methods.

As expected, real scenarios like those depicted in Distech-IR proved harder to generalize. A decrease in the performance was observed in all levels of supervision with a significant drop of 18.2% oLRP for *SSD* and 23.3% for *Yolo v5*. A smaller decrease was observed for AEs obtaining even closer results to the one from supervised approaches. The primary issue in this dataset was the large number of False Negative which almost doubled the FN obtained for FIR-Image-Action.

6 Conclusions

In this work, we presented a study comprising different methods with increasing levels of supervision for privacy-preserving person localization. Our experimental results over two low-resolution top-view IR datasets showed that reduced image-level supervision is enough for achieving results almost comparable to a fully-supervised detectors. Specifically, AE-based approaches proved to perform similarly to *Yolo v5* in real-world scenarios by only using images of empty rooms for training and with 3.7 times less execution time. Such a result is significant for reducing annotation costs and improving the scalability of intelligent building applications. Additionally, we detailed the process for producing bounding box annotations for low-resolution IR images and provided the localization for the publicly available dataset FIR-Image-Action.

Acknowledgements: This work was supported by Distech Controls Inc., and the Natural Sciences and Engineering Research Council of Canada (RGPIN-2018-04825).

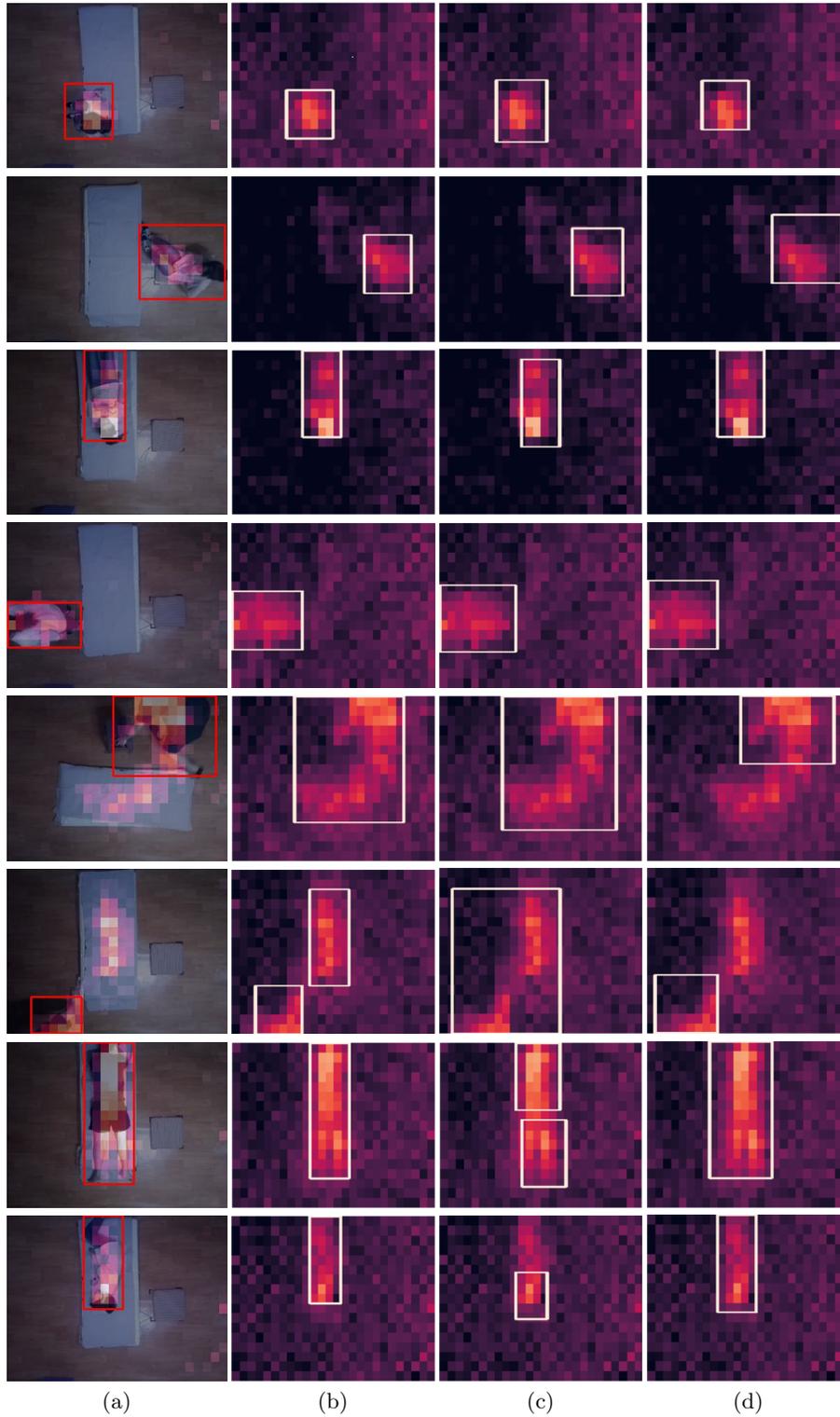


Fig. 5: Examples of low-resolution IR people detection results. Overlay of RGB and IR modalities with their corresponding ground truth (a), along with bounding box predictions of *dVAE* (b), *gradCAM* (c), and *Yolo v5* (d).

References

1. Bae, W., Noh, J., Kim, G.: Rethinking Class Activation Mapping for Weakly Supervised Object Localization. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (eds.) *Computer Vision – ECCV 2020*, vol. 12360, pp. 618–634. Springer International Publishing (2020)
2. Baur, C., Wiestler, B., Albarqouni, S., Navab, N.: Deep Autoencoding Models for Unsupervised Anomaly Segmentation in Brain MR Images. arXiv:1804.04488 [cs] **11383**, 161–169 (2019)
3. Belharbi, S., Sarraf, A., Pedersoli, M., Ayed, I.B., McCaffrey, L., Granger, E.: F-cam: Full resolution cam via guided parametric upscaling. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. pp. 3490–3499 (2022)
4. Buolamwini, J., Gebru, T.: Gender shades: Intersectional accuracy disparities in commercial gender classification. In: *Conference on fairness, accountability and transparency*. pp. 77–91. PMLR (2018)
5. Cao, J., Sun, L., Odoom, M.G., Luan, F., Song, X.: Counting people by using a single camera without calibration. In: *2016 Chinese Control and Decision Conference (CCDC)*. pp. 2048–2051. IEEE (May 2016)
6. Chen, J., Wu, J., Konrad, J., Ishwar, P.: Semi-coupled two-stream fusion convnets for action recognition at extremely low resolutions. In: *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*. pp. 139–147. IEEE (2017)
7. Chen, Z., Wang, Y., Liu, H.: Unobtrusive sensor-based occupancy facing direction detection and tracking using advanced machine learning algorithms. *IEEE Sensors Journal* **18**(15), 6360–6368 (2018)
8. Gao, C., Li, P., Zhang, Y., Liu, J., Wang, L.: People counting based on head detection combining adaboost and cnn in crowded surveillance environment. *Neurocomputing* **208**, 108–116 (2016)
9. He, P., Griffin, C., Kacprzyk, K., Joosen, A., Collyer, M., Shtedritski, A., Asano, Y.M.: Privacy-preserving object detection. arXiv preprint arXiv:2103.06587 (2021)
10. Jiang, P.T., Zhang, C.B., Hou, Q., Cheng, M.M., Wei, Y.: Layercam: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing* **30**, 5875–5888 (2021)
11. Jocher, G., Chaurasia, A., Stoken, A., Borovec, J., NanoCode012, Kwon, Y., TaoXie, Fang, J., imyhxy, Michael, K., Lorna, V, A., Montes, D., Nadar, J., Laughing, tkianai, yxNONG, Skalski, P., Wang, Z., Hogan, A., Fati, C., Mammana, L., AlexWang1900, Patel, D., Yiwei, D., You, F., Hajek, J., Diaconu, L., Minh, M.T.: ultralytics/yolov5: v6.1 - TensorRT, TensorFlow Edge TPU and OpenVINO Export and Inference (Feb 2022)
12. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
13. Kingma, D.P., Welling, M.: Auto-Encoding Variational Bayes (May 2014)
14. Kitano, H.: Classification and Localization of Disease with Bounding Boxes from Chest X-Ray Images p. 6
15. Lili Tao, Volonakis, T., Bo Tan, Ziqi Zhang, Yanguo Jing: 3D Convolutional Neural network for Home Monitoring using Low Resolution Thermal-sensor Array. In: *3rd IET International Conference on Technologies for Active and Assisted Living (TechAAL 2019)*. pp. 3 (6 pp.)–3 (6 pp.). Institution of Engineering and Technology, London, UK (2019)

16. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
17. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: SSD: Single Shot MultiBox Detector. arXiv:1512.02325 [cs] **9905**, 21–37 (2016)
18. Oksuz, K., Cam, B.C., Kalkan, S., Akbas, E.: One Metric to Measure them All: Localisation Recall Precision (LRP) for Evaluating Visual Detection Tasks. arXiv:2011.10772 [cs] (Nov 2021)
19. Otsu, N.: A threshold selection method from gray-level histograms. *IEEE Transactions on Systems, Man, and Cybernetics* **9**(1), 62–66 (1979)
20. Ren, Z., Lee, Y.J., Ryoo, M.S.: Learning to anonymize faces for privacy preserving action detection. In: Proceedings of the european conference on computer vision (ECCV). pp. 620–636 (2018)
21. Ryoo, M.S., Rothrock, B., Fleming, C., Yang, H.J.: Privacy-preserving human activity recognition from extreme low resolution. In: Thirty-First AAAI Conference on Artificial Intelligence (2017)
22. Schwemmer, C., Knight, C., Bello-Pardo, E.D., Oklobdzija, S., Schoonvelde, M., Lockhart, J.W.: Diagnosing gender bias in image recognition systems. *Socius* **6**, 2378023120967171 (2020)
23. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization. *International Journal of Computer Vision* **128**(2), 336–359 (Feb 2020)
24. Shengsheng Yu, Xiaoping Chen, Weiping Sun, Deping Xie: A robust method for detecting and counting people. In: 2008 International Conference on Audio, Language and Image Processing. pp. 1545–1549. IEEE, Shanghai, China (Jul 2008)
25. Sun, K., Zhao, Q., Zou, J.: A review of building occupancy measurement systems. *Energy and Buildings* **216**, 109965 (2020)
26. Tao, L., Volonakis, T., Tan, B., Jing, Y., Chetty, K., Smith, M.: Home Activity Monitoring using Low Resolution Infrared Sensor. arXiv:1811.05416 [cs] (Nov 2018)
27. Tateno, S., Meng, F., Qian, R., Hachiya, Y.: Privacy-Preserved Fall Detection Method with Three-Dimensional Convolutional Neural Network Using Low-Resolution Infrared Array Sensor. *Sensors* **20**(20), 5957 (Oct 2020)
28. Tateno, S., Meng, F., Qian, R., Li, T.: Human Motion Detection based on Low Resolution Infrared Array Sensor. In: 2020 59th Annual Conference of the Society of Instrument and Control Engineers of Japan (SICE). pp. 1016–1021. IEEE, Chiang Mai, Thailand (Sep 2020)
29. Wang, Z., Chang, S., Yang, Y., Liu, D., Huang, T.S.: Studying very low resolution recognition using deep networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4792–4800 (2016)
30. Yang, S., Kim, Y., Kim, Y., Kim, C.: Combinational Class Activation Maps for Weakly Supervised Object Localization. In: 2020 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 2930–2938. IEEE, Snowmass Village, CO, USA (Mar 2020)
31. Zhang, H.: Fir-image-action-dataset. <https://github.com/visiongo-kr/FIR-Image-Action-Dataset#fir-image-action-dataset> (2020)
32. Zheng, H., Zhong, X., Huang, W., Jiang, K., Liu, W., Wang, Z.: Visible-infrared person re-identification: A comprehensive survey and a new setting. *Electronics* **11** (2022)
33. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning Deep Features for Discriminative Localization. arXiv:1512.04150 [cs] (Dec 2015)