

Inter-Homines: Distance-Based Risk Estimation for Human Safety

Matteo Fabbri, Fabio Lanzi, Riccardo Gasparini, Simone Calderara, Lorenzo Baraldi, Rita Cucchiara

Abstract—In this document, we report our proposal for modeling the risk of possible contagiousity in a given area monitored by RGB cameras where people freely move and interact. Our system, called Inter-Homines, evaluates in real-time the contagion risk in a monitored area by analyzing video streams: it is able to locate people in 3D space, calculate interpersonal distances and predict risk levels by building dynamic maps of the monitored area. Inter-Homines works both indoor and outdoor, in public and private crowded areas. The software is applicable to already installed cameras or low-cost cameras on industrial PCs, equipped with an additional embedded edge-AI system for temporary measurements. From the AI-side, we exploit a robust pipeline for real-time people detection and localization in the ground plane by homographic transformation based on state-of-the-art computer vision algorithms; it is a combination of a people detector and a pose estimator. From the risk modeling side, we propose a parametric model for a spatio-temporal dynamic risk estimation, that, validated by epidemiologists, could be useful for safety monitoring the acceptance of social distancing prevention measures by predicting the risk level of the scene.

Index Terms—Computer Vision, Social Distancing, 3D people detection.

1 INTRODUCTION

THE COVID-19 emergency has changed the way we live interpersonal social relationships, at work, in public and private spaces, in places of education culture and leisure. The risk of contagion seems full-blown; until now, there are no conclusive studies which correlate environmental and endogenous factors with the greatest spread of the virus: instead, everything seems to correlate the contagion to proximity or to the contact between infected people and people susceptible to infection [1]. The spread of the infection seems to follow the epidemiological models that derive from the SIR models [2].

The phases that all the world is going to undertake after the lock-down will be characterized by living with the risk of contagion: the prerogative will be to take conscious and possibly interactive measures to minimise the possibility of contagion, while seeking a necessary resumption of social and working life.

Certainly, the IT technologies and in particular Artificial Intelligence can be valuable tools to monitor and predict risk levels in potentially crowded places. In fact, we propose an innovative and effective technological contribution based on Computer Vision and Deep Learning, in order to dynamic monitoring the acceptance of social distancing prevention measures through real-time calculation of the risk level, with particular reference to workplaces, public places and social areas. For statistical purposes, people behavior dynamics are stored in a database in a completely privacy compliant manner. The data can be used to identify the most critical areas and hours of the day in terms of number of people and risk level, in order to better address distance-related interpersonal prevention measures.

The system is called Inter-Homines (from the “Homo inter homines sum, capite aperto ambulo” - “I am a human among humans and I can walk with my face uncovered”) because people

should be free to move and interact with uncovered faces while being safe at the same time.

The system has a twofold goal. The first is to provide a reliable tool, in accordance with European privacy and usage guidelines of the AI, to calculate in real time the actual compliance with the prevention measures for “spacing”, also interactively reporting any risky situations. In particular, the implemented system can generate real-time alarms when people form crowds. The second goal is to provide an innovative model for the dynamic calculation of the risk of the monitored site that can be used as a tool for prevention, control, monitoring, and planning, support to the population and workers in order to implement conscious attendance, linked to effective compliance with the measures in force.

Detecting people, their position in the space, their mutual distance is a typical application of Computer Vision. Many tools are available, using state-of-the-art deep learning architecture and geometry-based 3D reconstruction. Results are promising although still far to be applied everywhere by everyone. In this project, we can take advantage of a long term experience in computer vision for surveillance and people behavior understanding [3], [4], providing a novel detection pipeline running in real-time. It exploits standard fast camera calibrations, a people detector and a pose estimation methods.

Inter-Homines defines a model, validated by epidemiologists and parameterizable according to current regulations, which allows, in real-time, to associate each monitored area with: a) a space-time risk index, b) a dynamic safety level of the area, c) a dynamic map of interpersonal distances and d) a real time visualization of detected persons and distances. See Fig. 1 for the system output overview.

2 RELATED WORK

One of the most popular two-stage deep object detectors is R-CNN [5] which predicts object location from a set of region proposals [6], crops them and classifies each using a second deep neural

• The authors are with the Department of Engineering “Enzo Ferrari”, University of Modena and Reggio Emilia, Italy.
E-mail: name.surname@unimore.it.

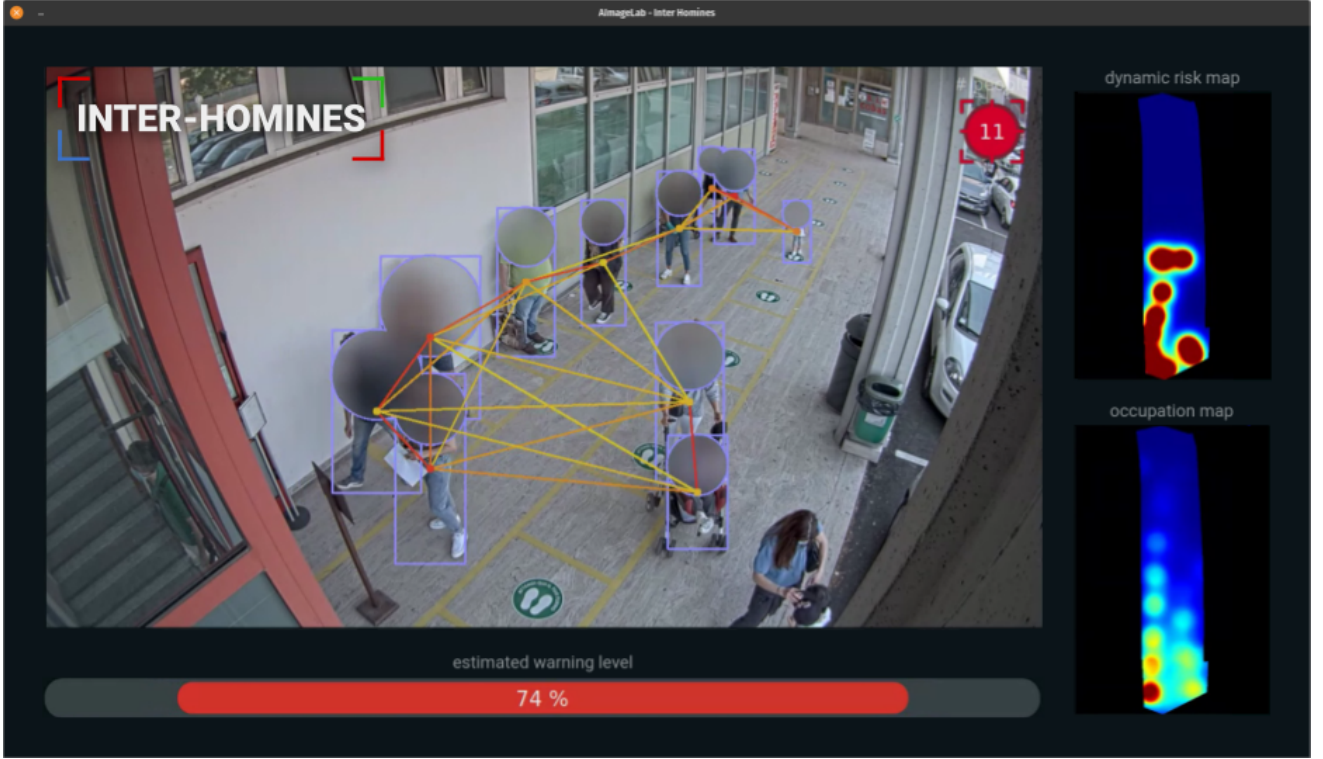


Fig. 1. GUI of our system. In the main frame, anonymized bounding boxes are superimposed to the image. Colored links encode people reciprocal distance. On the right, two maps show the bird-eye view of the area. The estimated risk level of the scene resides at the bottom of the interface.

network. Fast R-CNN [7], instead, directly crops image features to save computation. However, both approaches rely on slow low-level region proposal techniques.

On the other hand, one-stage methods such as Faster R-CNN [8] generates region candidates within the detection network. It samples bounding boxes with fixed shape (anchors) around the image grid and classifies them into foreground or background. Each proposal is then further classified into object classes. Several improvements to one-stage detectors include anchor shape priors such as in YOLO [9], [10], SSD's different feature resolution [11], and loss re-weighting among different samples [12].

Our approach leverages CenterNet [13], which is closely related to anchor-based one-stage detectors. However, CenterNet does not require manual thresholds for foreground and background classification and does not require Non-Maximum Suppression (NMS) [14] post processing as it simply extracts local peaks in the keypoint heatmap [3], [15]. Moreover, CenterNet utilizes an output stride of 4 which is 2 times larger than in traditional object detectors [16], [17], making it more accurate.

Other methods utilize the same robust keypoint estimation network as CenterNet: CornerNet [18] and ExtremeNet [19]. CornerNet detects the bounding box corners as keypoints while ExtremeNet predicts the left, top, right and bottom extremes of the objects. However, those methods require a combinatorial grouping stage as post processing, which considerably slows down the whole pipeline. CenterNet, instead, simply extracts a single center point per object without the need for grouping or post-processing.

People detection can be also achieved by pose estimation. The trend of pose estimation [15], [20], [21] is very promising, but often is too computational severe to be implemented for real time edge applications with an unknown number of people. Thus, in this work, we adopt a simplified pose estimation algorithm, that it

is used together with the people detector to make the localization more robust to occlusions.

Many 3D object detection methods have been proposed in literature. Among them, 3D R-CNN [22] adds a further head to Faster R-CNN [8] which is followed by a 3D projection. Also Deep Manta [23] exploits a coarse-to-fine Faster R-CNN [8] trained on multiple tasks. Finally, Deep3Dbox [24] utilizes a slow R-CNN [5] by first predicting 2D bounding boxes and then feeding each detection into a 3D estimation network. However, those methods require huge computational power and does not leverage constraints such as fixed camera and flat ground plane.

3 \mathcal{R}_0 AND THE SIR MODEL

After the outbreak of the COVID19 pandemic, all the world learned the importance of the basic reproduction number, \mathcal{R}_0 , as the statistical index indicating the degree of spread of the infection. In commonly used infection models, when $\mathcal{R}_0 > 1$ (in Italy has reached 4.3 during the spring of 2020) the infection will be able to start spreading in a population, but not if $\mathcal{R}_0 < 1$. Generally, the larger the value of \mathcal{R}_0 , the harder it is to control the epidemic.

\mathcal{R}_0 is defined as the expected number of secondary cases produced by an infection in a completely susceptible population:

$$\mathcal{R}_0 = \alpha \cdot c \cdot d \quad (1)$$

where α is the transmissibility, c is the average rate of contact between susceptible and infected individuals, and d is the duration of infectiousness.

To understand if this quantity defines the epidemic threshold of a particular infection, we need to formulate a Susceptible-Infected-Removed (SIR) epidemic model [25]. This model deploys several

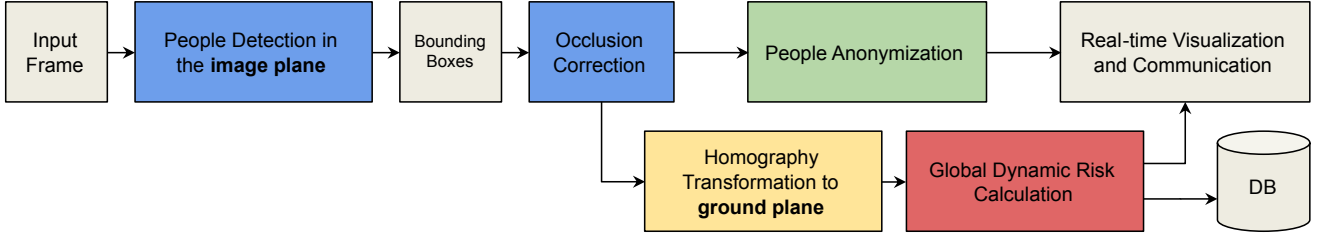


Fig. 2. Schematization of the Inter-Homines pipeline: the input frame is processed to produce bounding box detections. Each detection is then refined by the Occlusion Correction module that copes with truncated bounding boxes. The image plane detection coordinates are then transformed to ground plane coordinates using an Homography Transformation. Those coordinates are then used to calculate the global risk. People coordinates and risk level are then stored into a database. Finally, the system outputs the anonymized frame along with the risk level and the risk maps.

assumptions: 1) closed population, 2) constant rates, 3) no births and deaths and 4) well mixed population.

Given a population of N individuals, let's consider S the number of susceptible people, I the infected, and R the removed. Removed people are those that cannot be infected, as they might have developed antibodies. Now let's define $s = \frac{S}{N}$, $i = \frac{I}{N}$, $r = \frac{R}{N}$ as the fraction in each set. The SIR model is defined as:

$$\frac{ds}{dt} = -\beta si, \quad \frac{di}{dt} = \beta si - vi, \quad \frac{dr}{dt} = vi \quad (2)$$

An epidemic occurs if the number of infected increases: $\frac{di}{dt} > 0$. By considering that everyone is susceptible, we can substitute $s = 1$ obtaining the following inequality:

$$\alpha cd = \mathcal{R}_0 > 0 \quad (3)$$

\mathcal{R}_0 is essentially the entire theoretical basis of public health interventions for infectious diseases and it is simply the product of the transmissibility, the mean contact rate, and the duration of infection. In order to reduce transmissibility α we can develop vaccines, get people to use barrier contraceptives or use antivirals. To decrease mean contact c , the world decided to use isolation/quarantine, and health education programs. Finally, to reduce the duration of infection d therapeutics, antibiotic treatment of bacterial infections that boost innate immune response can be exploited.

\mathcal{R}_0 is in generally computed as a posterior measure, but cannot be dynamically predicted in a robust way since the factor influencing \mathcal{R}_0 are not a priori easily measurable. In this work we cannot do anything a part from monitoring the acceptance of health education programs. In the past months, many countries decided the mandatory measures of security that concern the use of DPI and the social distance guidelines. Thus, in order to make c as small as possible, we should keep all the people at a distance larger than a threshold distance of a possible infection.

A viable way is to force people to stay in queues, maybe with some marker placed on the floor and with the constant attention of a human guard that controls the compliance of the social distancing norms. This is not always possible, especially in big malls and wide areas. Moreover, the human monitoring is not always optimal as the guard is subject to tiredness and lost of focus.

This is the reason why computer based systems joined with risk models can substitute human controllers and help to perform real-time monitoring of areas by assessing a level of possible risk, and, if necessary, giving a real time feedback to improve the safety and decrease the risk. In the following section, we propose a very simple model, that, using some thresholds validated by epidemiologists, models the dynamic risk in a given area.

4 RISK MODEL

The SIR model formulation, as described in the previous section, has validity when considering a population. Now, let's consider a much more restricted zone. This could be an indoor area, like a waiting room of a public office, an entrance in a cinema or a shop. More precisely, let's consider a scene with N people k_0, \dots, k_{N-1} at a given time t . Given two subjects k_i and k_j with distance $d_{i,j}$, we define their reciprocal risk as follows:

$$rr_{i,j}^{(t)} = \eta e^{-\beta \max(0, d_{i,j} - \tau)} \quad (4)$$

where η , β and τ are parameters that respectively control height, slope and the full width at maximum of the function. In this specific application, η is a mitigator used to decrease the risk when some criteria are met, e.g., when at least one of the two people is wearing a facial mask. β , instead, controls how the risk decreases when the distance is greater than τ and can model environmental characteristics such as air temperature and the presence of air conditioning. Lastly, τ , controls the transmissibility of the disease via respiratory droplets and define the minimal distance allowed between two people. It should follow World Health Organization and national guidelines but can be further increased to better preserve the safety of people in critical places such as COVID-19 hospital units. We then define the individual risk at time t as:

$$R_i^{(t)} = \max_{j=0 \dots N-1, j \neq i} \{rr_{i,j}^{(t)}\} \quad (5)$$

The global risk at t of the scene is then computed as follows:

$$G^{(t)} = \min \left(1, \frac{1}{C} \sum_{i=0}^{N-1} R_i^{(t)} \right) \quad (6)$$

where C is the maximum capacity of the scene. This capacity can be either given by the user or calculated using simple covering algorithms. Finally, the dynamic global risk is computed as:

$$D^{(t)} = \frac{1}{W} \sum_{w=0}^{W-1} G^{(t-w)} \quad (7)$$

where W is the size of the temporal window. At a given time t , $D^{(t)} \in [0, 1]$ is the global risk of the scene and it is used to trigger alarms when it reaches a given threshold.

5 INTER-HOMINES TECHNICAL CORE

Here we give an overview of the pipeline we used to process videos in real time. The aim of our Inter-Homines system is to detect people, compute their distance and provide a dynamic risk level of the area, as well as producing a human readable visualization with anonymized people. For GDPR constraints,

no visual data is recorded but, instead, only people coordinates are extracted and stored. Data is acquired with a variable rate, up to one time per second for each camera. See Fig. 2 for a schematization of the pipeline.

The following subsections summarize the key elements of our system. Section 5.1 describes the people detection stage and elaborates on its challenges. Section 5.2 illustrates our proposed keypoint localization solution which addresses the occlusion problem peculiar of surveillance scenarios and also provide the head position for anonymization purposes. Next, in Section 5.3, we describe how we convert points from image plane to ground plane and, finally, Section 5.4 illustrates the system outputs.

5.1 People Detection

As we are interested in the best speed-accuracy trade-off, we choose CenterNet [13] as a people detector. In particular, we rely on the DLA backbone [26] which yields 51.3% AP for the people class on MS COCO [27], running at 52 FPS on a Titan XP.

Let $I \in \mathbb{R}^{W \times H \times 3}$ be the input image having width W and height H . CenterNet outputs a keypoint heatmap $\hat{Y} \in [0, 1]^{\frac{W}{R} \times \frac{H}{R} \times C}$, where $R = 4$ is the output stride and C is the number of keypoint types. Keypoint types include $C = 80$ object categories but in this work we only consider the “people” class. Detected keypoints corresponds to a prediction $\hat{Y}_{x,y,c} = 1$ and 0 otherwise. To recover the discretization error generated by the output stride, CenterNet further predicts a local offset $\hat{O} \in \mathbb{R}^{\frac{W}{R} \times \frac{H}{R} \times 2}$ for each center point.

Let $(x_1^{(k)}, y_1^{(k)}, x_2^{(k)}, y_2^{(k)})$ be the bounding box of object k of the “people” class and $p_k = (\frac{x_1^{(k)} + x_2^{(k)}}{2}, \frac{y_1^{(k)} + y_2^{(k)}}{2})$ it’s center point. CenterNet predicts all center points for each object k and further regresses to the object size $s_k = (x_2^{(k)} - x_1^{(k)}, y_2^{(k)} - y_1^{(k)})$.

At running time, we first extract the peaks in the heatmap for the “people” category. We consider all the responses whose value is greater or equal to its 8-connected neighbors. Let $\hat{\mathcal{P}} = \{(\hat{x}_i, \hat{y}_i)\}_{i=1}^n$ be the set of n detected center points where keypoint values $\hat{Y}_{x_i y_i c}$ are utilized as a measure of its detection confidence. Bounding boxes are produced at location:

$$\begin{aligned} (\hat{x}_i + \delta\hat{x}_i - \hat{w}_i/2, \hat{y}_i + \delta\hat{y}_i - \hat{h}_i/2, \\ \hat{x}_i + \delta\hat{x}_i + \hat{w}_i/2, \hat{y}_i + \delta\hat{y}_i + \hat{h}_i/2), \end{aligned} \quad (8)$$

where $(\delta\hat{x}_i, \delta\hat{y}_i) = \hat{O}_{\hat{x}_i, \hat{y}_i}$ is the predicted offset and $(\hat{w}_i, \hat{h}_i) = \hat{S}_{\hat{x}_i, \hat{y}_i}$ is the predicted size. Since the prediction are directly produced from the keypoint estimation, there is no need for IoU-based NMS or other post-processing techniques. This makes CenterNet faster w.r.t. other detectors, making it suitable for real time applications.

CenterNet is capable of producing a precise localization of every person in the image, however, it does not take into account occlusions that usually happen in real world scenarios. As shown in Fig. 3 (pink bounding boxes), if a person is occluded by an object or by other people, CenterNet predicts a tight bounding box that only contains the visible part of the person, ignoring his full shape. This usually happens with the bottom part of the body, as the camera is commonly placed several meters above the ground. Since we are ultimately interested in recovering the ground plane coordinate of each person through homography, we need to know the exact position (in image plane) of the feet of each detected person. This task cannot be accomplished by solely relying on CenterNet.



Fig. 3. Examples of CenterNet bounding boxes (pink), refined bounding boxes and head localization (green).

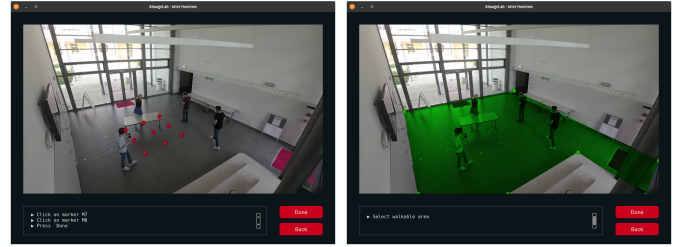


Fig. 4. GUI used during system calibration for homography matrix calculation (left) and for walking area selection (right).

5.2 Feet and Head Localization

To overcome the aforementioned limitations without introducing complexity to the overall system, we propose to utilize a small network to predict the feet position given a bounding box containing a person, even if the feet are not visible.

To this aim we rely on a simple but effective CNN that, given an image M tightly containing a person, it regresses to the midpoint $P_f = (x_f, y_f)$ of the segment having the two feet as endpoints. This ensures that we know the exact position in image plane where every person touches the ground. Since we are also interested in anonymizing the face of each detected person, we further predict the location of the head $P_h = (x_h, y_h)$.

We replaced the last 1000 class classification layer of Resnet50 [16] with two heads composed by an adaptive average pooling layer and a fully connected layer with output size equal to 2. The adaptive average pooling takes care of the difference in size that each bounding box fed to the network can have. Training has been carried out for 10 epochs using an MSE loss with Adam optimizer, batch size of 64 and learning rate 0.001.

We used JTA [3] as the training dataset since it is the only surveillance dataset available in literature that provide pose estimation annotations with occlusion information. Thanks to this, we were able to simulate occlusion situations by simply picking, during training, the pedestrians with the bottom keypoints occluded, like ankles, knees, and hips. During training, we also randomly shortened some of the bounding boxes in order to simulate CenterNet behaviours. This step ensures a more precise localization of the feet while also coping with truncated bounding boxes. As shown in Fig. 3 (green bounding boxes), our network can effectively obtain an accurate position of each head and it is used to extend the bounding box to its regular shape.

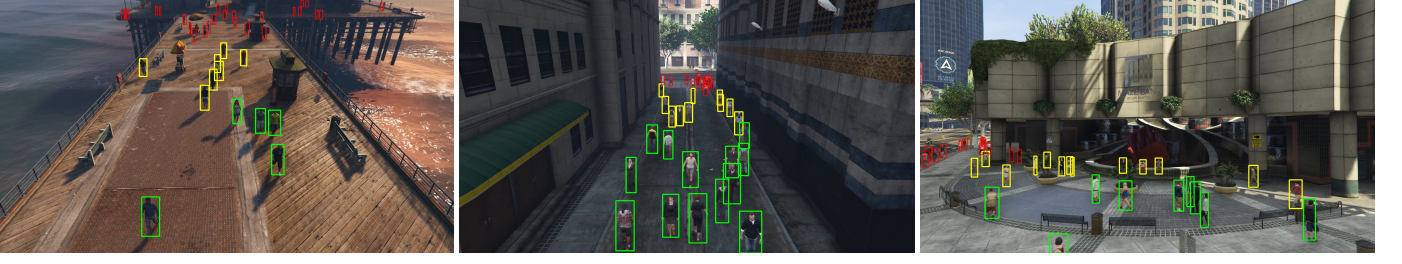


Fig. 5. Examples from the JTA dataset exhibiting its variety in viewpoints, number of people and scenarios. Ground truth bounding boxes are superimposed to the original images. Green color is used for people having a distance from the camera between 0 and 20 meters, yellow for people between 20 and 40 meters and red for people between 40 and 100 meters.

5.3 From Image Plane to Ground Plane

The camera projection matrix P is a 3×4 matrix which describes the mapping of a pinhole camera [28] from 3D points in the world to 2D points in an image. Let X be a representation of a 3D point in homogeneous coordinates (a 4-dimensional vector), and let y be a representation of the image of this point in the pinhole camera (a 3-dimensional vector), we have $y = PX$. The camera projection matrix can be decomposed as:

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} r_{1,1} & r_{1,2} & r_{1,3} & t_X \\ r_{2,1} & r_{2,2} & r_{2,3} & t_Y \\ r_{3,1} & r_{3,2} & r_{3,3} & t_Z \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad (9)$$

where the intrinsic parameters f_x , f_y and c_x , c_y are the camera focal length and principal points respectively while $r_{i,j}$ and t_i are the extrinsic parameters which define the rotation and the translation used to describe the rigid motion of an object in front of a still camera. Finally, u and v are the coordinates of the projected point in pixels while X , Y and Z are the coordinates of a 3D point in the world coordinate space. By considering the simpler case of a projection of planar points, where each 3D point lies on the same plane (e.g. the ground), we can simplify the formulation considering $Z = 0$. For planar surfaces, 3D to 2D perspective projection reduces to a 2D to 2D transformation:

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{1,1} & r_{1,2} & t_1 \\ r_{2,1} & r_{2,2} & t_2 \\ r_{3,1} & r_{3,2} & t_3 \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} \quad (10)$$

and by doing the products we finally obtain the planar homography matrix H . The planar homography relates the transformation between two planes (e.g. the image plane and the ground plane):

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} h_{1,1} & h_{1,2} & h_{1,3} \\ h_{2,1} & h_{2,2} & h_{2,3} \\ h_{3,1} & h_{3,2} & h_{3,3} \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} = H \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} \quad (11)$$

Since H maps from ground plane to image plane, but we are interested in the opposite transformation (from image plane to ground plane), we now need to calculate the inverse homography matrix H^{-1} . An homography matrix H is always invertible, and its inverse is still an homography transformation:

$$\begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} = H^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \quad (12)$$

A practical way to calculate the homography matrix H of Eq. 11 is to find a set of at least four points pairs of target and source planes and to minimize the back-projection error:

$$\sum_{i=0}^N \left[\left(u_i - \frac{h_{1,1}X_i + h_{1,2}Y_i + h_{1,3}}{h_{3,1}X_i + h_{3,2}Y_i + h_{3,3}} \right)^2 + \left(v_i - \frac{h_{2,1}X_i + h_{2,2}Y_i + h_{2,3}}{h_{3,1}X_i + h_{3,2}Y_i + h_{3,3}} \right)^2 \right] \quad (13)$$

However, if not all of the point pairs fit the rigid perspective transformation, this initial estimate will be poor. To solve this problem we employ the RANSAC iterative method, trying many different random subsets of the corresponding point pairs (of four pairs each). We then estimate the homography matrix applying a simple least-square algorithm using this subset, and then compute the quality of the computed homography, which is the number of inliers. The best subset is then used to produce the initial estimate of the homography matrix. The computed homography matrix is refined further (using only the inliers) with the Levenberg-Marquardt method to further reduce the re-projection error. The homography matrix is determined up to a scale. Thus, it is normalized so that $h_{3,3} = 1$.

This method of using an homography transformation to obtain 3D coordinates is the most appropriate when we want to monitor an approximately flat area (such as a town square) using a fixed camera and there is the possibility of making appropriate measurements in the monitored space.

To easily obtain the points pairs of image and ground planes needed to find the homography matrix H , we designed a simple procedure that we call “system calibration”. This procedure consists in placing nine markers at the center of the monitored area, fully visible from the camera. The markers are placed in a grid pattern as in Fig. 4. By means of a simple graphic interface, the user can take a snapshot of the camera and click with the cursor the centers of the nine markers in order to acquire the pixel coordinates. In practice, we utilize a special carpet with the nine markers printed on it. The use of the carpet automates the real world measurements as we already know the distance between markers in the carpet, making the system calibration fast, less prone to errors and feasible by everyone. Once the nine pairs of points have been identified and the homography matrix calculated, the carpet can be safely removed and the system will continue to work properly as long as the camera maintains its position.

During the system calibration an optional procedure of selecting the “walking area” can be carried out. Again, a simple graphic interface let the user draw a polygon on the snapshot

TABLE 1

3D detection results on JTA Dataset. In PR (precision), RE (recall) and F1, @ t indicates that a predicted person is considered “true positive” if the distance from the corresponding ground truth location is less than t . The max range indicates the maximum distance considered in the calculation.

		PR	RE	F1	PR	RE	F1	PR	RE	F1
max range		@0.5 m			@1.0 m			@1.5 m		
10m	Inter-Homines w/o Occ. Corr.	83.38	78.29	80.06	90.44	85.28	87.07	92.56	87.63	89.31
	Inter-Homines Full Pipeline	88.01	84.74	85.55	92.95	89.2	90.22	94.27	90.56	91.59
20m	Inter-Homines w/o Occ. Corr.	69.11	59.75	63.41	85.33	73.86	78.36	91.83	79.77	84.48
	Inter-Homines Full Pipeline	74.96	66.98	70.00	88.86	79.20	82.87	93.64	83.69	87.49
30m	Inter-Homines w/o Occ. Corr.	59.47	46.87	51.57	77.39	60.76	66.96	85.81	67.27	74.14
	Inter-Homines Full Pipeline	65.3	53.07	57.52	81.26	65.4	71.18	88.34	70.96	77.31
100m	Inter-Homines w/o Occ. Corr.	53.76	31.65	38.07	71.34	41.74	50.41	80.3	46.88	56.7
	Inter-Homines Full Pipeline	60.91	36.21	43.37	77.12	45.21	54.55	84.75	49.41	59.79

taken from the camera, as shown in Fig. 4. The pixel vertices are then converted to ground coordinates that are used to exclude detections whose 3D position is outside the walking area. This is useful, for example, to ignore mirrors or windows that can reflect people causing unwanted detections.

Given a bounding box of a person, we can now extract its central point (u, v) of the lower side of the box (i.e. the image coordinate where the person touches the ground with his feet), and utilize Eq 12 to obtain the corresponding (X, Y) point in ground plane. Now that we have the 3D position of every person in the scene, the dynamic global risk in Eq. 7 can be calculated and given as output along with other information that we summarize in the following subsection.

5.4 System Output

A convenient graphical interface highlights all the main results of the analysis of our Inter-Homines system, allowing to evaluate at a glance the crowding conditions in the monitored area (see Fig. 1). This interface is made with Qt to guarantee compatibility with all the main operating systems.

5.4.0.1 Anonymized Frame: It shows real-time the bounding box detections superimposed to the input RGB frame. The system is privacy compliant and all the faces are obscured. Colored segments connect people who are at an estimated distance lower than a defined upper threshold distance (typically 3 m). The color indicates the extent of the infraction, going from a dark red for the most serious infraction to yellow for the minor ones.

5.4.0.2 People Counter: At the top right of the frame we also display the number of detected people updated in real time. This number is an average computed in a window of W frames to account for miss detections and false positives.

5.4.0.3 Dynamic Risk and Occupation Maps: In the right part of the interface two bird eye views of the walking area are updated real-time. The Dynamic Risk map shows a snapshot of the current situation of the area. The Occupation map, instead, displays the overall aggregated risk and it is computed by averaging the Dynamic Risk maps of the whole day. It is used to identify areas with a larger risk for statistical and predictive purposes. Note that people outside of the walking area are completely ignored and do not affect the statistics.

5.4.0.4 Estimated Warning Level: In the lower part of the window a bar represents the total estimated risk in the monitored area and it is computed using Eq. 7. The application provides the possibility to send an alarm signal (example: send an email / audio

notification) if certain thresholds on the number of people or on the risk level are exceeded. The thresholds and the notification methods of their exceeding will be defined according to the needs of the context in which the system will operate.

5.4.0.5 Weekly Report: Since we want to give insightful statistics to help with the prevention of the infection, our system periodically produce a report. The report contains statistics about number of people, risk level, number of infractions and occupation maps aggregated by hours and days. To this end we utilize a non-relational database to store timestamp and position of each person captured by our system.

6 SYSTEM VALIDATION

In order to validate the effectiveness of our system, we performed a series of experiments leveraging JTA [3]. JTA is a massive dataset for pedestrian pose estimation and tracking in urban scenarios created exploiting the highly photorealistic video game *Grand Theft Auto V*. The videos feature a vast number of different people appearances, in several urban scenarios at varying illumination conditions and viewpoints. Each clip comes with a precise annotation of visible and occluded body parts, people tracking with 2D coordinates in image plane and 3D coordinates in camera space. JTA overcomes all the limitation of existing datasets in terms of number of entities and available annotations. Each video contains a number of people ranging between 0 and 60 with an average of more than 21 people, totaling almost 10M annotated body poses over 460,800 densely annotated frames. The distance from the camera ranges between 0.1 and 100 meters, resulting in pedestrian heights between 20 and 1100 pixels. JTA is composed by a set of 512 Full HD videos, 30 seconds long, recorded at 30 fps.

As shown in Fig. 5, despite being a synthetic dataset, JTA features highly challenging and complex situations, peculiar of surveillance scenarios, where people are often dominated by severe body part occlusions and truncations. We believe this dataset is the perfect choice to validate a system that targets global safety.

Since we can not perform the system calibration procedure on an already recorded dataset, i.e. we can not physically place the markers at the center of the scene, we designed a simple heuristic to directly recover the nine points pairs using the dataset annotations. With the assumption that every foot of each person lies on the same plane, for each JTA sequence, we linearly regressed the ground plane utilizing the 3D coordinates of every foot in every frame of that sequence. By recovering a unit normal vector of the plane and two orthonormal vectors lying on the plane

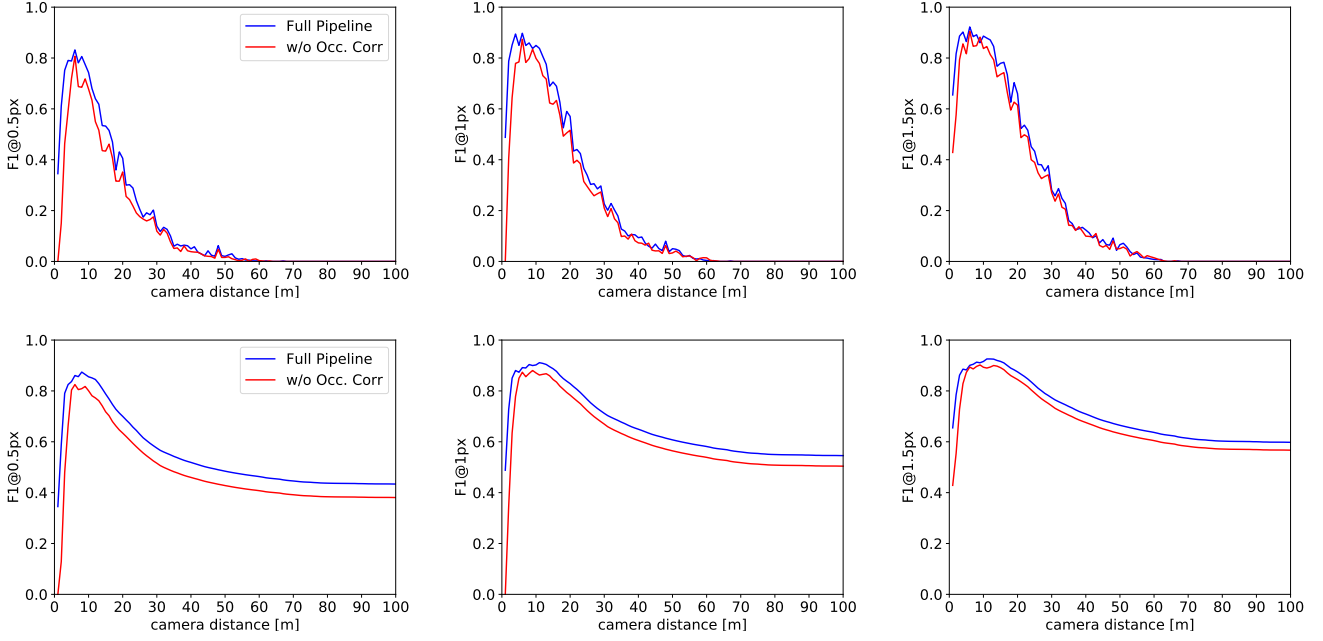


Fig. 6. F1 score vs. camera distance at different thresholds (first row) and F1 score vs. max camera distance at different thresholds (second row).

we were able to find the orthonormal base of the new space that allowed us to move each 3D coordinate into a space where each foot has the same y coordinate (according to the standard camera system). Now, since each foot coordinate has the same y , we can get rid of it and considering the new (x, z) coordinates as ground coordinates. As we are interested in nine points pairs of target and source planes, we utilized a K-Means implementation to find nine foot cluster centers. Utilizing a clustering method ensures that the nine points are far from each other. Once recovered the foot cluster centers, we remapped those coordinates into the original standard camera space and projected them into the image plane using the pinhole camera model. The 2D projected coordinates and the 2D foot clusters now form the nine points pairs needed to calculate the homography matrix.

Experiments are conducted on every 10th frame of a subset of the JTA test set where we carefully removed the sequences that contain camera motion and people at different heights, e.g. people going up the stairs, as our method assumes static camera and flat ground plane. Tab. 1 shows the precision, recall and F1 obtained using different thresholds and considering different camera distance ranges. As the range increase, we observe a decrease in performances, due to the fact that small people are hardly detected and homography transformation becomes less reliable. Since we are interested in evaluating the impact of that occlusions have in the performance of our system, we reported the results with and without the occlusion correction module. As can be shown, the correction is always beneficial, especially when people are close to the camera.

To better understand how performance degrades as distance increases, in Fig. 6, first row, we plotted the F1 score at different thresholds w.r.t. the camera distance. It is interesting to note that performance worsens when people are too close to the camera. In Fig. 6, second row, we plotted the same quantity but, this time, the F1 score is calculated considering all the people with distance less than the camera distance, and not equal to the camera distance.

7 CONCLUSIONS

In this work we proposed a simple and effective system that deals with the COVID-19 emergency by providing a social distancing tool that can prevent the spread of the infection. We validated it using a highly challenging benchmark, obtaining a lower bound on the performance of the method. We believe that our system can be a practical solution to an important problem, hoping to see areas less crowded than the JTA dataset in the near future.

ACKNOWLEDGMENTS

The work is supported by the Italian MUR, Ministry of Universities and Research, under the project PRIN 2019-2021 “PRE-VUE Prediction of Events in Urban Environment” project and partially supported by EU European Regional Development Funds for Regione Emilia Romagna under the special project “Inter-Homines” 2020 among the 16 research and innovation projects for the development of solutions aimed at contrasting the epidemic of COVID-19.

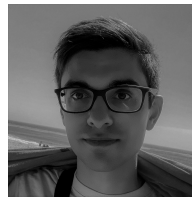
REFERENCES

- [1] S. Asadi, N. Bouvier, A. S. Wexler, and W. D. Ristenpart, “The coronavirus pandemic and aerosols: Does covid-19 transmit via expiratory particles?” 2020.
- [2] Y.-C. Chen, P.-E. Lu, C.-S. Chang, and T.-H. Liu, “A time-dependent sir model for covid-19 with undetectable infected persons,” *arXiv:2003.00122*, 2020.
- [3] M. Fabbri, F. Lanzi, S. Calderara, A. Palazzi, R. Vezzani, and R. Cucchiara, “Learning to detect and track visible and occluded body joints in a virtual world,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [4] M. Fabbri, F. Lanzi, S. Calderara, S. Alletto, and R. Cucchiara, “Compressed volumetric heatmaps for multi-person 3d pose estimation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (ECCV)*, 2020.
- [5] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014.

- [6] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Subcategory-aware convolutional neural networks for object proposals and detection," in *IEEE winter conference on applications of computer vision (WACV)*, 2017.
- [7] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015.
- [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*, 2015.
- [9] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [10] —, "Yolov3: An incremental improvement," *arXiv:1804.02767*, 2018.
- [11] 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*, 2016.
- [12] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017.
- [13] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv:1904.07850*, 2019.
- [14] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis, "Soft-nms—improving object detection with one line of code," in *Proceedings of the IEEE international conference on computer vision*, 2017.
- [15] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [17] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *IEEE international conference on computer vision*, 2017.
- [18] H. Law and J. Deng, "Cornersnet: Detecting objects as paired keypoints," in *Proceedings of the European Conference on Computer Vision*, 2018.
- [19] X. Zhou, J. Zhuo, and P. Krahenbuhl, "Bottom-up object detection by grouping extreme and center points," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [20] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *European conference on computer vision*, 2016.
- [21] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, "Deepcrut: A deeper, stronger, and faster multi-person pose estimation model," in *European Conference on Computer Vision*, 2016.
- [22] A. Kundu, Y. Li, and J. M. Rehg, "3d-rcnn: Instance-level 3d object reconstruction via render-and-compare," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- [23] F. Chabot, M. Chaouch, J. Rabarisoa, C. Teuliere, and T. Chateau, "Deep manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [24] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka, "3d bounding box estimation using deep learning and geometry," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [25] W. O. Kermack and A. G. McKendrick, "A contribution to the mathematical theory of epidemics," *The royal society of london. Series A, Containing papers of a mathematical and physical character*, 1927.
- [26] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, "Deep layer aggregation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- [27] 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*, 2014.
- [28] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.



Matteo Fabbri is currently a Ph.D. student at the International Doctorate School in ICT of the University of Modena e Reggio Emilia. He works under the supervision of Prof. Rita Cucchiara and Prof. Simone Calderara, on computer vision and deep learning for people behavior understanding. He worked 10 months at Panasonic Silicon Valley Lab as a Deep Learning Engineer. His research interests include generative models, pose estimation and multiple object tracking.



Fabio Lanzi is currently a research fellow at the Artificial Intelligence Research and Innovation Center (AIRI) promoted by the "Enzo Ferrari" Department of Engineering and by the "Marco Biagi" Department of Economics of the University of Modena and Reggio Emilia. Since 2017, he works under the supervision of Prof. Rita Cucchiara and Prof. Simone Calderara, and he has he took part in a series of industrial research projects mainly concerning human pose estimation, tracking and action recognition.



Riccardo Gasparini is currently a research fellow at Almagelab, a research laboratory of the Department of Engineering "Enzo Ferrari" at the University of Modena and Reggio Emilia, Italy. He works under the supervision of Prof. Rita Cucchiara on various topics such People Detection, Tracking, People Recognition, Video Surveillance, Anomaly Detection, Video Analysis, Egocentric vision and Embedded sensors.



Simone Calderara received a computer engineering masters degree in 2005 and the Ph.D. degree in 2009 from the University of Modena and Reggio Emilia, where he is currently an assistant professor within the AlmageLab group. His current research interests include computer vision and machine learning applied to human behavior analysis, visual tracking in crowded scenarios, and time series analysis for forensic applications. He is a member of the IEEE.



interests include image processing, video understanding, deep learning and multimedia.

Lorenzo Baraldi received the Ph.D. degree (cum laude) in information and communication technologies from the Università degli studi di Modena e Reggio Emilia, Italy, in 2018. He was a Research Intern with Facebook AI Research (FAIR) in 2017. He is currently an Assistant Professor with the Dipartimento di Ingegneria Enzo Ferrari, Università degli Studi di Modena e Reggio Emilia. He has authored or coauthored over 50 publications in scientific journals and international conference proceedings. His research



The research carried out spans on different application fields, such as video surveillance, automotive and multimedia big data annotation. Currently, she is AE of IEEE Transactions on Multimedia and serves in the Governing Board of IAPR and in the Advisory Board of the CVF.

Rita Cucchiara received the masters degree in Electronic Engineering and the Ph.D. degree in Computer Engineering from the University of Bologna, Italy, in 1989 and 1992, respectively. Since 2005, she is a full professor at the University of Modena and Reggio Emilia, Italy, where she heads the AlmageLab group and is Director of the CINI AIIS Lab. She published more than 300 papers on pattern recognition computer vision and multimedia, and in particular in human analysis, HBU, and egocentric-vision.