



C³Net: End-to-End Deep Learning for Efficient Real-Time Visual Active Camera Control

Christos Kyrkou

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

May 21, 2022

C³NET: END-TO-END DEEP LEARNING FOR EFFICIENT REAL-TIME VISUAL ACTIVE CAMERA CONTROL

THIS PAPER IS A PREPRINT OF A PAPER SUBMITTED TO AND ACCEPTED FOR PUBLICATION IN JOURNAL OF
REAL-TIME IMAGE PROCESSING.

Christos Kyrkou*

KIOS Research and Innovation Center of Excellence
Department of Electrical and Computer Engineering
University of Cyprus
1 Panepistimiou Avenue, Nicosia Cyprus
{kyrkou.christos}@ucy.ac.cy

July 28, 2021

ABSTRACT

The need for automated real-time visual systems in applications such as smart camera surveillance, smart environments, and drones necessitates the improvement of methods for visual active monitoring and control. Traditionally, the active monitoring task has been handled through a pipeline of modules such as detection, filtering, and control. However, such methods are difficult to jointly optimize and tune their various parameters for real-time processing in resource constraint systems. In this paper a deep Convolutional Camera Controller Neural Network is proposed to go directly from visual information to camera movement in order to provide an efficient solution to the active vision problem. It is trained end-to-end without bounding box annotations to control a camera and follow multiple targets from raw pixel values. Evaluation through both a simulation framework and real experimental setup, indicate that the proposed solution is robust to varying conditions and able to achieve better monitoring performance than traditional approaches both in terms of number of targets monitored as well as in effective monitoring time. The advantage of the proposed approach is that it is computationally less demanding and can run at over 10 FPS ($\sim 4\times$ speedup) on an embedded smart camera providing a practical and affordable solution to real-time active monitoring.

Keywords Real-Time Active Vision, Smart Camera, Deep Learning, End-to-End Learning

1 Introduction

Real-time active vision systems (i.e., movable cameras with controllable parameters such as pan and tilt) can provide extended coverage, flexibility, and cost-efficiency compared to static vision systems running in the cloud. Active cameras can track targets (i.e. follow them) in order to record their movements and alert in case of an intrusion [1]. Real-time operation on resource-constraint hardware is needed for such systems which are increasingly being used for various applications ranging from surveillance [2], mobile robots [3], and intelligent interactive environments [4].

Currently existing approaches for active vision decompose the problem into separate modules, namely detection, tracking, and control and employ different algorithms for the purpose of detecting targets and then following them [5, 6]. Such examples include motion detection, background modelling/subtraction, and lastly tracking by detection. The former two are widely used in static camera scenarios [7] due to their relative computational efficiency however, for active cameras that move and exhibit constant background change the latter methods are preferred since they are more widely applicable. Such techniques are often augmented with post-processing filters that increase the computational complexity making it difficult to use in embedded scenarios where the algorithm needs to run on the camera itself in real-time [8]. Hence, existing methods that rely on hand-crafted features, motion models, and modeling of the camera

*ckyrkou@gmail.com, www.christoskyrkou.com

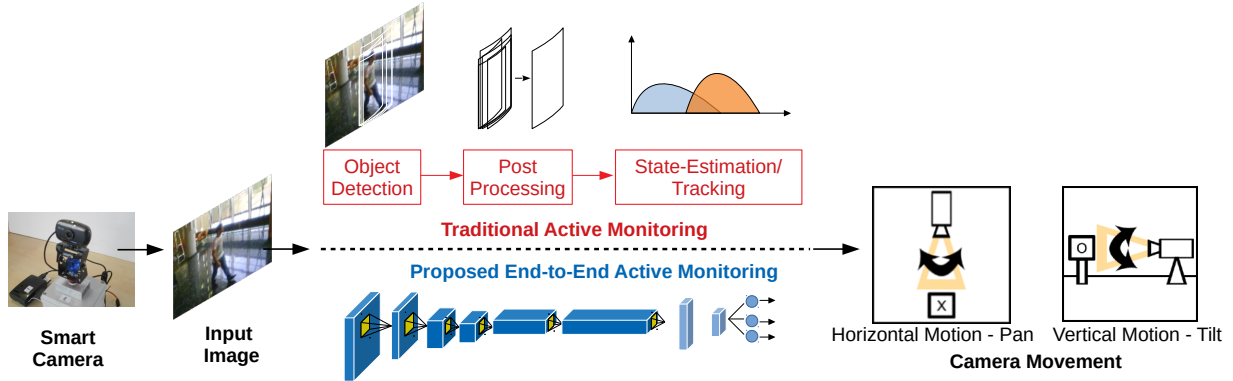


Figure 1: Traditional active monitoring Vs Proposed end-to-end deep visual active camera control.

views are not optimized for active vision scenarios e.g., tracking in the case of a camera mounted on a pan-tilt station [7]. Some approaches utilize a master-slave paradigm, but there has not yet been any attempt to deal with active monitoring of multiple targets in an end-to-end way. End-to-end learning approaches allow encapsulating all the intelligence into the machine learning algorithm thus can optimize all processing steps simultaneously and learning the best features to associate with camera control for visual active monitoring purposes.

To deal with the aforementioned challenges, this paper investigates the appealing potential of end-to-end learning using Convolutional Neural Networks (CNN) [9, 10, 11] to develop an active vision system for monitoring people in surveillance and smart-environment applications. In particular the task of visual active monitoring is formulated as a supervised learning problem to be solved by deep learning in an end-to-end manner so that the low-level image features are directly associated with camera motion. In doing so it is possible to have a more compact system with reduced sub-modules to optimize and more efficient operation. An efficient Convolutional Camera Controller neural network, referred to as C^3 Net, is designed and trained in an end-to-end manner to associate image features with control actions. To achieve the former an appropriate simulation framework and dataset for training and evaluating CNN-based controllers for active smart cameras is also introduced.

The effectiveness of the proposed C^3 Net has been verified by experiments through the simulation framework using image sequences from a publicly available PETS2009 [12] dataset. Moreover, experiments were also performed under a real setup of a smart camera implementation and UAV acquired data. Results indicate that the proposed end-to-end approach is able to consistently monitor 1 – 3 targets more on average than other methods and for a longer period of time. In addition, through experiments designed to test the robustness of the method it is validated that the network learns to focus on the majority of targets, without been given the bounding boxes and never explicitly trained to detect people. It also facilitates real-time performance for surveillance applications using resource-constraint platforms as it is computationally efficient as it has a small size and operates at over 10 frames-per-second (FPS).

The rest of this paper is structured as follows. Section II outlines relevant research and key areas of background material. Section III formally introduces the problem as well as the data collection and training procedures, and the proposed controller CNN. The experimental results and evaluation are presented in Section IV. Section V provides a discussion on possible extensions, while Section VI concludes the paper.

2 Background and Related Work

2.1 Problem Overview

The objective of active vision in contrast to a static camera setting is to change the control parameters of a vision sensor in order to maximize a visual-task-related performance objective such as following one or more targets that are located within the field-of-view (FoV) to improve the overall surveillance and monitoring capabilities. In particular this work deals with *active monitoring* defined as the process of controlling the camera in order to position the center of its Field-of-View (FoV) at the center of mass of one or more targets in the image [13] (i.e., keep subject(s) as close to the center of the camera image as possible).

Typically, detection algorithms [14] that provide bounding box estimates centering around objects (Figure Fig. 1) are used to guide the active vision systems. Bounding boxes used for localization and estimating target position can be

affected by noise, occlusions, and even bad Non-maximum-suppression parameters [14]. In addition, hand-crafted post-processing steps such as non-maximum suppression, used to remove superfluous detections, can lead to localization errors that affect the active monitoring performance [15, 16]. Furthermore, tracking and state estimation approaches require fine tuning various parameters and can be negatively affected by the detection performance degradation. Therefore, there is a need to improve visual active monitoring methods for emerging applications that simultaneously require real-time performance, have battery limitations or operate in remote locations and with requirements for rapid deployment in temporary installations by providing computational cost that can lead to simpler active vision systems.

2.2 Visual Active Monitoring

Over the last years, deep neural networks, especially Convolutional Neural Networks (CNN), have improved the state-of-the-art in static object tracking/monitoring [17, 18, 19, 6]. Conventional solutions for active visual tracking tackle the problem by decomposing it into two or more sub-tasks [20], i.e., object detection typically using a machine-learning-based classifier/detector, a tracking algorithms such as Kalman filter [21], and a control output for the camera movement. Each task is optimized individually resulting in highly complex systems with many tuning parameters. Furthermore, this leads to difficulty in obtaining real-time performance on resource-constraint embedded camera systems. Different works have investigated the use of active cameras with one or more degrees of freedom for monitoring applications and are summarized next.

Initial approaches such as [22], followed a master-slave approach to track targets at a high resolution. One camera, the master, has a wide FoV and performs blob detection and uses a Kalman filter to track a target in an area and by projecting from image plane to world coordinates it controls (pan, tilt, and zoom) the other active camera, the slave, to follow a target. In contrast the goal of this work is to improve the tracking performance of a single camera agent so that it can autonomously follow all targets present in its FoV.

In [23], the authors proposed a camera control module to follow an intruder. A background model is used to detect motion in the scene while a Proportional-Integral (PI) controller moves the camera to follow the target. Illumination changes and crowded environments can affect the performance of such techniques. Moreover, the approach is only suitable for a single target.

A similar system was developed in [24] where a combination of Haar cascades and color histograms are used to track face and upper body of a target in indoor environments. A fuzzy controller is used to steer the camera. The input to the fuzzy controller are the x and y position, as well as the size of the target object in the image. Likewise, the output of the fuzzy controller are the required pan, tilt, and zoom settings for the camera. The behavior of the system is determined by a set of rules connecting input values to expected outputs. The system can only track a single target and the hand-engineered rules can hinder its performance.

The system proposed in [25] uses face detection to find a single face in the camera. It then extracts features for tracking and then feedback control loop is applied to vary the pan and tilt camera angles by certain steps based on the target velocity until it is centered.

The approach in [26] uses complementary features and a multi-scale mechanism to follow the subject in front of the camera. In the first stage, a discriminative pre-filter automatically chooses the features which best separate object and background to build a confidence map which helps to remove all low confidence samples. In the second stage, a rough model of object appearance is built at the beginning and continuously updated according to all changes in pose, viewpoint, illumination. The accumulation over time can cause the model to fail with sudden changes in pose and cluttered background. Furthermore, this system again tracks only a single object.

In [27] an Extended Kalman filter is used to jointly track the object position in the real world as well as estimating the intrinsic camera parameters. The filter outputs are used as inputs to two PID controllers (one for pan and one for tilt motion axis) which continuously track a moving target at a certain resolution. The focus of this work is on tracking only a single target however, and the objects are assumed to have a predetermined size.

In [5] the authors propose a control approach for visual tracking to estimate and reduce the impact of different disturbances affecting the servoing system performance. The experiments were conducted with a pan-tilt camera following a single moving target. Furthermore, there is an assumption that the tracking module always performs ideally.

The active camera system proposed in [4] is composed of multiple components in order to track a subject. The face of the target is identified through a face detection system and then a tracker is employed to estimate the targets motion across frames based on previous observations. An online learning approach is also used to learn the appearance model over time and reduce false detections. In addition, Gaussian Mixture Models are used to model the body movement in case the face detection fails. The final part is the controller of the pan-tilt camera which makes its decision of how to

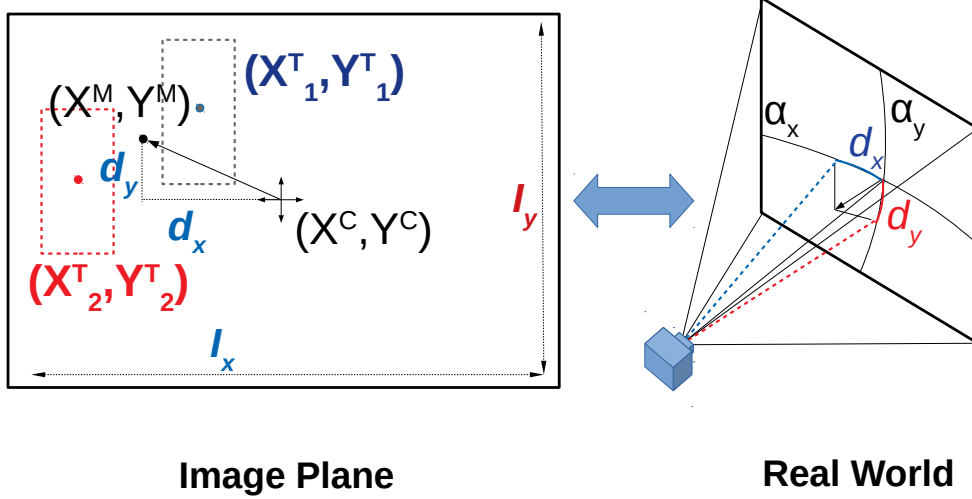


Figure 2: Mapping from Image Plane to Real-World: (Left) Displacement between image center and target center of mass. (Right) Camera angle of views associated with motion at each axis.

move based on where in the image the target is positioned and using a set of rules. This multi-component system can be difficult to tune and transfer to more resource-constraint systems that monitor multiple targets.

An end-to-end active tracker is proposed in [28] for following a character in the VizDoom video game environment. It uses reinforcement learning to train a CNN with LSTM to output discrete movement actions. Even though they are somewhat realistic, these scenarios do not correspond to real-world use cases since evaluation is done in a video game environment whereas herein images of real environments from PETS2009 are used. Furthermore, the step-based output controls a video game player and is not suitable for the dynamic range need to control a pan-tilt camera. In addition, the control actions happen within the context and dynamics of a virtual world. In contrast, in this work the active monitoring problem is formulated in such a way as to be directly applicable to real-world camera systems.

In summary, it is evident from the literature that related works make excessive use of multiple modules composed of hand-crafted models and rules that must be tuned separately and in most cases track only a single target [6]. While there has been considerable progression in utilizing deep learning for static camera tracking there has been relatively few works dealing with deep learning for active smart camera systems. In addition, most of them do not consider important system requirements, such as real-time constraints on low-end hardware systems.

This work attempts to bridge the gap between the use of active cameras with deep learning algorithms by proposing an end-to-end learning approach to simultaneously built an implicit detector and controller for cameras with pan and tilt motion capabilities. In contrast to existing works e.g. [28, 6] that follow a single target our goal is to give the same priority to all targets and attempt to follow as many as possible without a specific focus a single target. In addition, it does not require to be given an anchor target to follow and thus can be used in generic scenarios where the goal is to monitor targets in an area. Overall, the approach eliminates the need to rely on bounding box predictions which can be noisy and face difficulties with dense targets.

3 Deep Active Visual Monitoring

3.1 Approach Overview

The problem of end-to-end active visual control of a camera through a single image, as shown in Fig. 1, is formulated as a regression problem. The input is an image I from the camera sensor with resolution $I_x \times I_y$, and the output is a motion control vector \vec{M} . The control vector corresponds to the pan and tilt motion that the camera will have to perform in order to position the target(s) close to the center of the image. A Convolutional Camera Controller neural network (C^3 Net) is trained end-to-end to learn a control function f such that $\vec{M} = f(I)$.

A pinhole camera model and rectilinear lens [29] are assumed, as illustrated in Fig. 2. Under this assumption the pixel distances of the image center to the camera center are analogous to the angle that the servo motors have to move in order to position the target at the center of FoV [5, 25, 3, 30]. In addition, the motion of the camera when following targets in

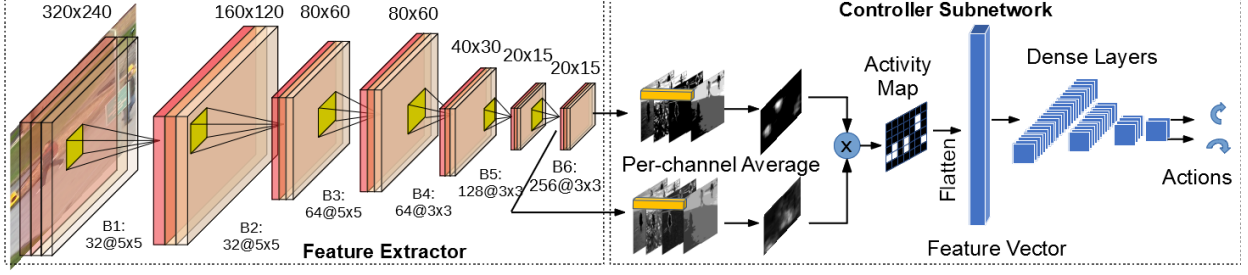


Figure 3: C³Net Architecture: Conceptually comprised of a feature extractor and a controller that condenses the high level semantic information to infer an action.

its FoV is bounded by the half of its horizontal and vertical angles of view (α_x, α_y). Accordingly, the motion of the camera in the the pan and tilt axes (denoted by M_x and M_y) can be calculated with respect to the viewing angles [30]. This can be done by first calculating the distance of the camera center from the target(s) center of mass (X^M, Y^M) for the horizontal and vertical direction (d_x, d_y), and then use Eq. 1 to associate the pixel distance to the camera angles. Hence, during the learning process the objective is to map the input image to normalized pixel displacement values ($d_x/I_x, d_y/I_y$) which will then be used to calculate the corresponding angle displacements in the pan and tilt axes. Predicting offsets instead of angles effectively bounds the output making it easier for the network to learn, and decouples the learning process from camera specific parameters.

$$M_x \rightarrow \frac{d_x}{I_x} \times \frac{\alpha_x}{2}, M_y \rightarrow \frac{d_y}{I_y} \times \frac{\alpha_y}{2} \quad (1)$$

3.2 Convolutional Camera Controller Network (C³Net)

C³Net is a suitable CNN trained to perform the regression task and estimate the camera motion displacement. Figure 3 shows the network architecture which is comprised of two main parts, a feature extractor made up of convolutional layers and a motion controller that summarizes the feature maps in order to calculate the final output values. Since the system is trained end-to-end it is difficult to distinguish the parts of the network that function primarily as feature extractor and which serve as controller. The input to the CNN is an RGB image which has normalized pixel values between $[0 \dots 1]$. The network can process images of 320×240 , but any other image size can be provide after it is resized.

Convolutional Feature Extractor: The layers were designed to perform feature extraction and were chosen empirically through a series of experiments that varied layer configurations. There are 5 major blocks each comprised of a convolutional layer with *Relu* activation and batch normalization layer. Furthermore, dropout is applied at the middle of the sub-network to combat overfitting with a rate of 0.2. Overall, the feature extractor is designed to be inherently computationally efficient to support use in embedded smart cameras for local control and decision making. To reduce the computational cost the image is down-sampled early on and a relatively small number of filters is used in the convolutional layers. Note that this is not an explicit object detector that produces bounding boxes, but rather implicitly learns the features to focus on through the controller supervision.

Controller Subnetwork: The controller subnetwork is responsible to map the extracted image features to control actions. It is comprised of layers with predetermined functionality as well as fully-connected layers. The custom layers serve to condense the information from the feature extractor essentially encoding our prior knowledge on the problem. The custom layers perform a per channel average of the feature maps extracted from blocks 5 and 6 (Fig. 3). The reason of using the last two layers is that they encode more semantically meaningful information thus can provide better indication for the presence of objects in the image. The response of each feature map is averaged over the channels and then multiplied element-wise emphasizing the overlapping regions. This activity map is then reshaped into a vector and given as input to the fully connected layers. The idea behind this is to leverage multi-level features to first encode an activity map with the presence of objects over an image grid and then provide a small vector to the fully connected layers to reduce parameter count. There are 4 fully connected layers with 100,50,10 and 2 neurons respectively to map the feature vector to motion controls. There is a dropout in between the dense layers and all layers have a *leakyrelu* activation function. The output of the controller subnetwork is further processed through a hyperbolic tangent activation that bounds the output between $[-1, \dots, 1]$. Overall, C³Net has a total of $\sim 500,000$ parameters which requires ~ 4 MB, resulting in a lightweight network that can run even on low-end CPUs.

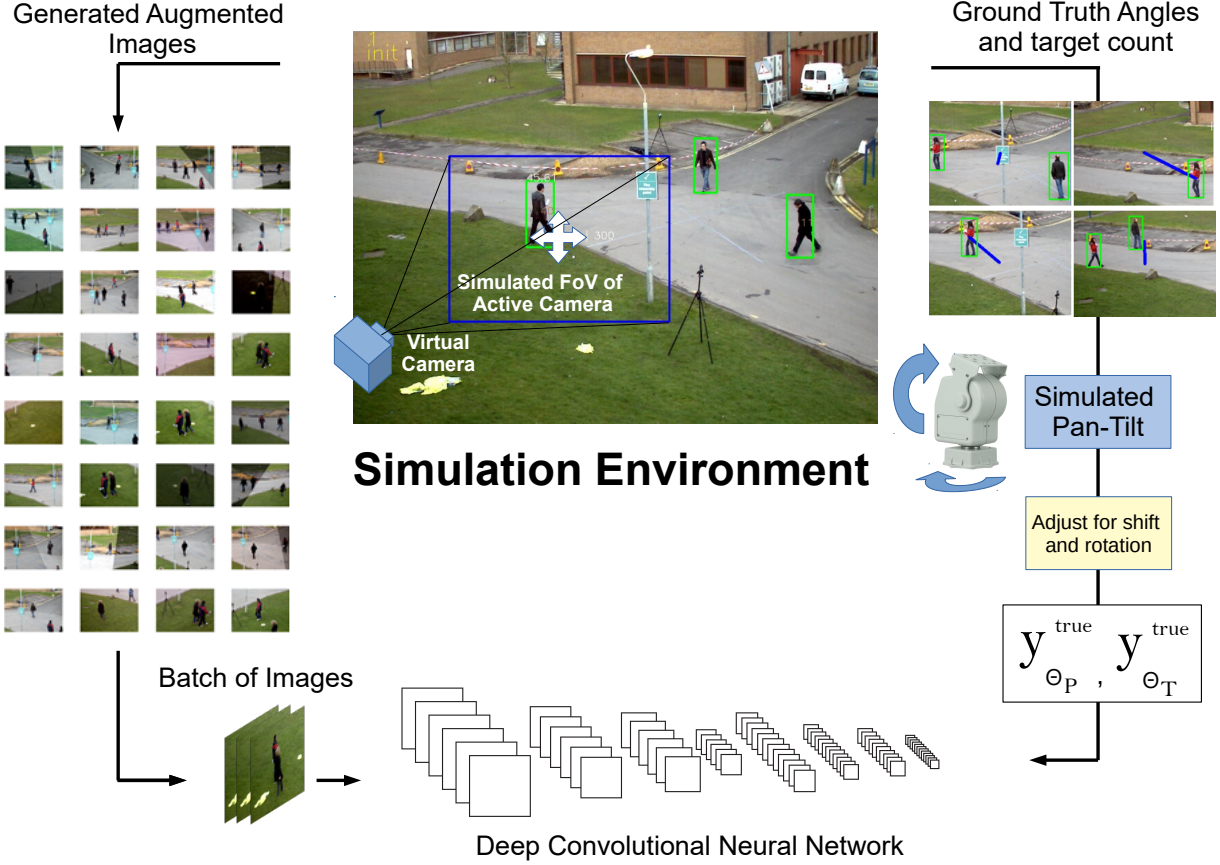


Figure 4: Overall framework for simulating virtual active camera movements and tracking, and extracting ground truth data.

3.3 Active Camera Data Generation

Proper training and testing data are necessary to train a deep CNN regressor for the visual active monitoring task. To the best of the authors knowledge there is no publicly available dataset for active vision applications with ground truth camera controls. For this reason an existing multi-person image dataset, such as the PETS2009[12], traditionally used for static tracking [31] is re-purposed to develop a simulation framework. The particular dataset is selected over more recent ones as it has multiple targets, larger image frames where targets only occupy a small region, and has ground truth bounding box annotations.

Overall, the framework allows for i) simulate the behaviour of active cameras using real-world images, ii) Capturing and storing multiple frame sequences with ground truth data that can be used for bounding box, density, and camera control, iii) evaluate performance of active vision algorithms in a realistic environment in similar conditions and controlled experiments. The simulation framework shown in Fig. 4 simulates the movement of a virtual active camera with pan-tilt control with a fixed FoV by moving it within a larger image frame. By cropping a fixed region from the original frame, it effectively restricts the active camera to a limited part of the overall image frame. These crops take $\sim 15\%$ of the original image. Obtaining the ground truth data for each image crop is a straight forward procedure if bounding box annotations are provided with the dataset, as is the case with PETS2009. First, the center of mass is calculated for all objects within each crop. The difference between the centre of the mass of the annotated people and the centre of the cropped region, in x, y directions, dictates the two control signals for that frame as per Eq. 1. Alternatively if no bounding boxes are provided it is possible to use a trained object detector to provide candidate boxes and train in a weakly supervised manner. The ground truth of images without objects is set to $(0, 0)$, meaning no movement. It is worth noting the same exact values used to move the camera in the simulated environment can be used to move the camera in the real-world when multiplied with the camera horizontal and vertical viewing angles. Thus the framework provides a valuable complementary data source for developing and benchmarking data-driven methodologies for smart camera applications.

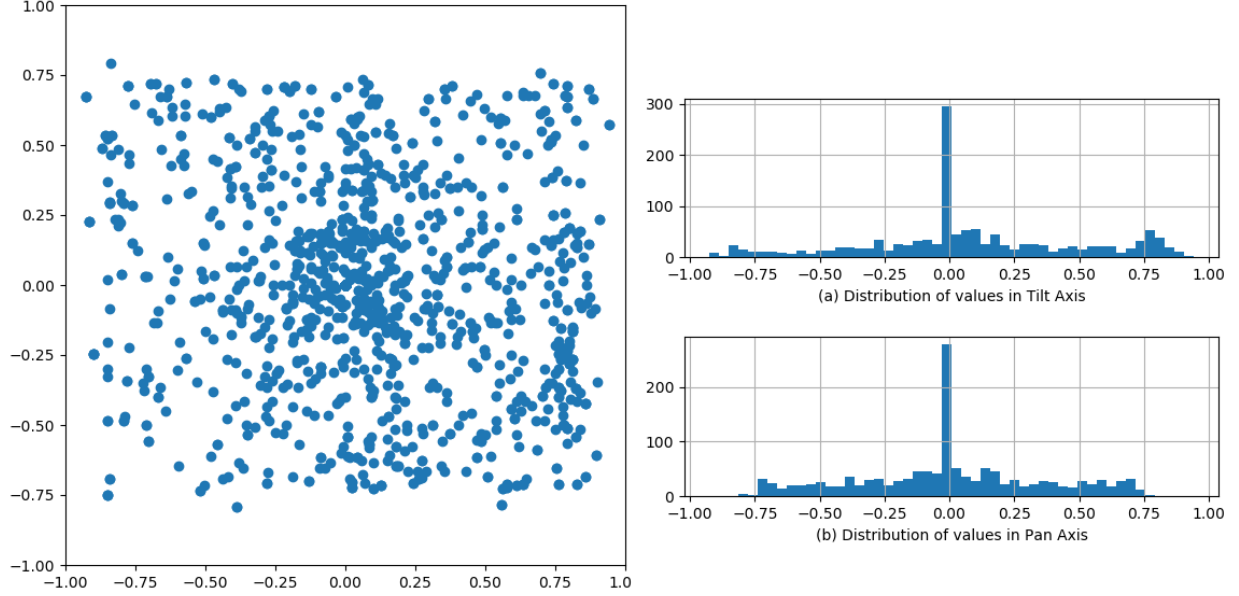


Figure 5: Distribution of the pan and tilt values. The majority of values are close to zero so selective sampling is used to balance the dataset.

A single sequence from the PETS-2009 [12] database, with ground truth bounding boxes, was selected as the basis to generate training data. A dataset of over 4500 frames was generated using the simulation software in order to train C^3 Net. Starting at various positions in random frames in the sequence the changes in displacement of pan and tilt axes are calculated according to Eq. 1 so that the targets will be positioned at the camera image center in the next frame. Hence, the training set effectively contains single images sampled from the sequence, paired with the corresponding ground truth future pan and tilt displacement. The pan and tilt displacement values are normalized between -1 and 1 with regards to the maximum motion of the camera in both axes. For this work images of 320×240 are considered, which have been used in similar studies [24], but it is also trivial to extract different resolution images.

An important part of data preparation is balancing the data. During the data collection process depending on the motion of successive frames, the target displacement values may not change significantly in both axes. As a result, the distribution of camera motion values might be imbalanced as shown in Fig. 5. To ensure that the learning algorithm does not overfit by producing small values around zero, during training the data in each batch is sampled to contain high as well as low values.

The whole data is randomly split into 60% for training the CNN regression model, and 40% of the data was used for validation purposes. In addition, augmentations are probabilistically applied to the images to increase the variability and combat over-fitting. The augmentation strategy included some transformations on the image pixels such as blurring and sharpening, color-shifting, illumination changes. Geometric transformations such as translations and horizontal flip were also performed with appropriate adjustment of the target pan and tilt values. The combinations of all these augmentations resulted in a variety of novel images that were used for training.

3.4 Network Training

The objective of the learning process is to regress a motion vector \vec{M} corresponding to motion in the pan and tilt axes. Accordingly the euclidean distance loss function in Eq. 2 is employed for learning the camera controls between ground truth y and predictions \hat{y} . The *Keras* deep learning framework [32] with *Tensorflow* [33] running as the backend is used for the training of the CNN regressor. The network was trained using a GeForce Titan Xp, on a PC with an Intel *i7 - 7700K* processor, and 32GB of RAM. The Adam optimization method was used for training with a learning rate step-decay approach starting from an initial learning rate of 0.001, and decreasing it by 0.5 when encountering a plateau. The CNN regressor is trained for 300 epochs with a batch size of 128 and 50 batches per epoch resulting in 6400 augmented images per epoch.

$$L = \frac{1}{N_B} \sum_{j=1}^{N_B} \left[(y_{d_x}^j - \hat{y}_{d_x}^j)^2 + (y_{d_y}^j - \hat{y}_{d_y}^j)^2 \right]^{1/2} \quad (2)$$

3.5 Weighted Moving Average Filtering

A form of momentum is introduced to smooth out the outputs of the neural network and avoid being misled by a single frame, that is based on a weighted moving average (Eq. 3). It is grounded on the assumption that the target motion will not change rapidly between frames. By incorporating a simple form of memory into the system the previous actions are also considered which reduces the impact of possible outliers. A window of 3 most recent outputs is maintained ($K = 3$) which was empirically found to represent a good trade-off between having a smoother control output but not being influenced too much by passed control outputs. The weight w can be determined through different policies. Herein, the weight factor is determined as shown in Eq. 3, where t equals K is the most recent output that is weighted the highest. Each frame the network produces output \hat{y} , and the moving average produces the filtered output \bar{y} .

$$\bar{y} = \sum_{t=1}^K \hat{y}_t \times w_t, \text{ where } w_t = \frac{t}{\sum_{z=1}^K z} \quad (3)$$

4 Evaluation

To compare the proposed method a traditional active tracking pipeline is implemented as shown in Fig. 1. It follows the tracking-by-detection approach used by related works for active camera control [34, 21]. In this paradigm an object detector is used to localize the targets and a tracker is applied to filter the detections over time. Finally based on the filter target positions the camera is moved to position the targets to their center of mass. We compare the different approaches with respect to their monitoring efficiency with metrics outlined in section . Also through an experimental setup we compare the performance of each approach in terms of real-time performance on an embedded smart camera.

The remaining sequences from the PETS2009 dataset that were not used during training and have available ground truth data are used to generate a separate testing dataset. The sequences differ in the time of day, illumination conditions, density and amount of targets, paths of the targets as well as the identity, scale, and appearance of the targets. As such, the performance of the proposed approach is measured objectively. In addition, additional datasets were collected using two different experimental setups to further test whether the approach generalizes to completely different environments. The proposed approach is compared against three other baselines. All have a detector and tracking component. Comparisons are shown with since single-shot detectors the smaller YOLO[35] variant referred to as *tinyYOLO*. In addition, for fairness a CNN (referred to as *PDN*) for pedestrian detection is trained with the YOLO framework on the same dataset as C^3 Net and is developed from scratch for detecting pedestrians on 320×240 images. It has a slightly larger feature extraction subnetwork and the latter stages are different to regress the bounding boxes compared to C^3 Net. Comparisons are also made against the widely used SVM-HOG pedestrian [36] detector implemented in OpenCV [37] still used in many applications (e.g., [38, 39]). On top of the detectors tracking is applied with Kalman filter which is a common approach used in active tracking [29, 27, 21]. The incorporated filtering handles bounding box associations, maintains trajectories, and handles the creation and termination of tracks. The kalman filter is preferred since it was not only use in prior works for active tracking but also is less computationally and memory demanding while providing comparable performance [21]. Also note that there is no explicit detection step implemented within the proposed network. As such, it is not possible to directly compare detection performance with standard object detection approaches and metrics. Furthermore, our goal is not detection itself but rather the monitoring performance. It is important to note that we restrict the comparison to methods followed by previous works and that are more suitable for smart camera applications where the compute-budget is limited. It is worth noting that recent deep-learning-based methods for tracking running on GPU platforms achieve 2 – 5 FPS [19]. Hence, methods such as visual trackers that rely on cascade CNNs [40] that increase the processing time on top of the detection part are not considered. Another important differentiation to recent works that utilize CNNs for tracking is that in their majority they rely on static camera images and specifically on tracking single objects in the frame. Hence, they are not directly comparable.

4.1 Performance Metrics

The image frames are taken from the PETS2009 dataset however, performance is not measured in the same way since the objective is not to compare with traditional tracking systems that operate on static sensors. Specifically, besides

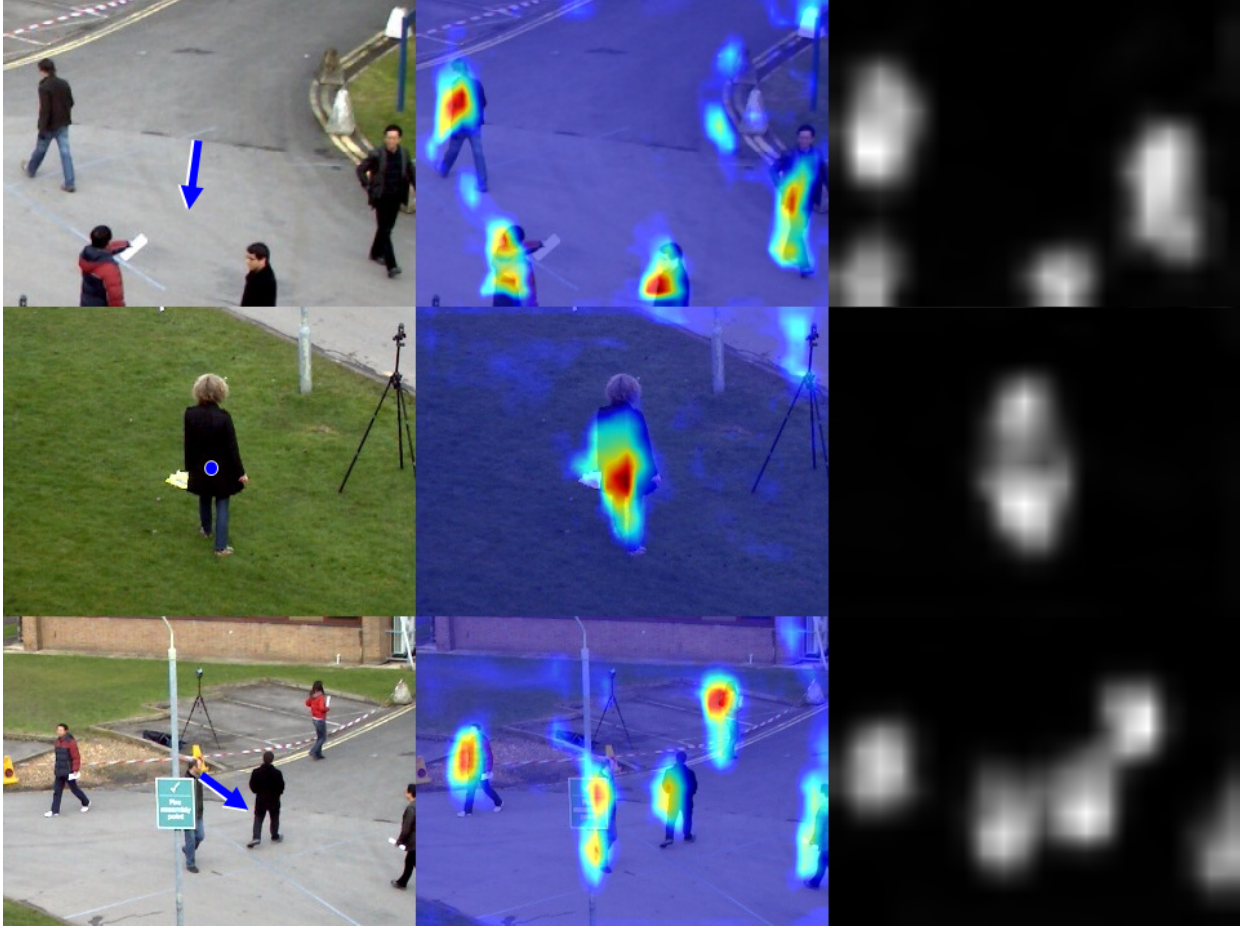


Figure 6: Visualizing what the network learns to look for in images (left) Input Image and direction of motion estimated by C^3 Net. (Center) Visualizations created using Grad-CAM for finding the pixels that mostly influence the network prediction. (Right) Activation map produced by the controller sub-network by aggregating the features from B5 and B6.(Best viewed in color)

calculating the error between ground truth and predicted motion vector from a set of test images, three metrics are also examined that can holistically evaluate how well an algorithm manages to monitor targets within an area. For all metrics for a target to be considered visible over 50% of the target’s body should appear in the frame

1) **Average Number of Targets in FoV:** The mean number of targets that are monitored over the duration of the sequence. It is the average across the number of visible targets at each frame. It essentially measures how well an algorithm can handle multiple targets within a frame.

2) **Average Monitoring Time:** The percentage of time for which one or more targets are visible within the FoV of the camera. It essentially measures how well an algorithm does in keeping up with the target(s). This is important aspect as it provides an indication of how likely it is to lose a target and how likely it is to follow a new one once it has entered the FoV.

3) **Average Distance from Target Centroid:** This is the Euclidean distance of the image center (i.e., camera center) to the centroid of the visible targets (i.e., center of mass). It shows how well an algorithm can keep the targets at the center of its FoV.

4.2 Evaluation of Learning

First, a closer look is taken into what features and image regions influence the prediction of C^3 Net and what it has learned to respond to in the various cases. To this end, the Gradient-weighted Class Activation Mapping (Grad-CAM) approach proposed in [41] is employed to generate heat maps of the image regions that mostly influence the control



Figure 7: Results on the same image with and without the presence of targets. In the case that no target is present in the image the camera does not move. (Best Viewed in Color)

decision by the network. The visualization is with regards to the target outputs of the network. In addition, the activation maps of the feature-extracting sub-network are visualized to understand what features the CNN is looking at (specifically the sum of the spatial averages from blocks B5 and B6). As shown in Fig. 6, the Grad-CAM visualization the network bases its control output on the targets in the image even if they are partially occluded. Note, that the network is never given the bounding boxes and never explicitly trained to detect people or any other body part but as shown in the activation visualization these are exactly the features that it detects in order to control the camera. Another way to gain confidence in what the network has learned to look for is through an ablation study where image parts are systematically removed to see how relevant they are to the network output. As shown in Fig. 7 for images not used during training, the network will only steer the camera to move only when a target is present in the image view. Also notice that it responds even in cases where the target is at different resolutions and not fully visible.

4.3 Performance in Simulated Environment

The next experiment involved evaluating all the aforementioned methods in the simulation environment and the performance metrics from Sec. 4.1. But first it is important to understand why common tracking benchmarks are not well suited for this purpose. The main objective during evaluation is to examine the behaviour of each algorithm and the influence of the control actions on the overall performance. Hence, a camera view needs to be placed within a larger frame in a suitable simulation environment that provides images of moving targets that the virtual camera can follow. In addition, the set of images needs to come with bounding box annotations to calculate whether a target is within the FoV and subsequently measure the different performance metrics described in section 4.1. Hence to conduct these experiments, 3 sequences from the PETS2009[12] dataset were used for which the ground truth is available but have not been used in the training and validation phases. The developed simulation framework outlined in section 3.3 models the camera motion based on the visual input and evaluates its performance. The output of each vision pipeline is a motion vector that will be passed to the simulator to perform the action. These sequences provide an additional challenge as they feature different visual conditions with higher crowd densities and different motion patterns. All approaches start at the same point in every video sequence. The images in the sequences are of 768×576 resolution and the virtual camera FoV is set to 320×240 for all methods so that there is margin for the camera to move and follow the targets. The first sequence (Seq. 1) has 241 images and features sparse groups of people walking together, the second (Seq. 2) has 795 frames and shows individuals independently moving in the area, and the third (Seq. 3) contains 436 frames of sparse crowd moving randomly in different directions. The camera FoV for all methods is set at the same initial position. In all cases the objective of the camera is to move in such a way as to keep the most number of targets in its FoV. Fig. 8 shows some examples from each sequence and the direction of movement chosen by the network.

Performance on Still Images: First to understand precisely how well each algorithm performs, comparisons are made against the ground truth motion on a per image basis and the overall error for each motion direction is calculated. Table 1 shows that the proposed CNN achieves the lowest average error as well as the lower maximum error in terms of estimating the camera motion controls. Since still images are used the detection approaches are mostly influenced by the predicted bounding boxes, their size, and the smoothing of the response maps through the non-maximum suppression process. In the case of C^3 Net the only difference is that the moving average is not employed. Overall, as shown in

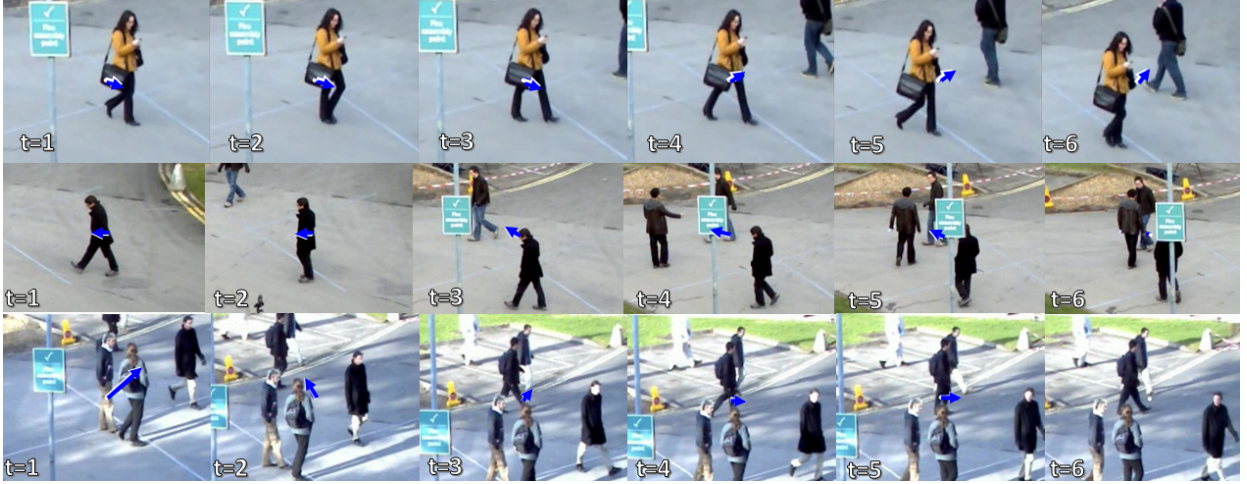


Figure 8: Samples of Sequences: (a) Seq1: Following 1 and then 2 targets. Notice that the network decides to move towards the center of mass of the two targets. (b) Seq2: The network follows a single target then it moves to the center of the group. Also the network does not move away from the partially occluded target. (c) Seq3: Multiple targets present the network adjusts to the movement of incoming targets.

Table 1: Error analysis on estimated values.

Method	Avg. Error		Max Error	
	θ_x	θ_y	θ_x	θ_y
C ³ Net (Proposed)	0.074	0.020	0.315	0.344
SVM-HOG w tracking	0.170	0.086	0.917	0.784
tinyYOLO w tracking	0.136	0.077	0.899	1.076
PDN w tracking	0.076	0.034	0.459	0.351

Table 1, the proposed approach has a much lower error rate across the two control signals while in the worst case it still has lower error margins.

Performance on Continuous Frames: Results using the metrics defined in 4.1 on continuous frames from the 3 aforementioned sequences are shown in Table 2 which will also incorporate the temporal aspect and demonstrate the effect of how compounded errors can affect the overall performance. First, the average number targets in the FoV is recorded which indicates how well each approach manages to follow a group of targets. The proposed approach manages to outperform the other approaches by following between 1-5 more targets on average as it is not affected by missed bounding box detections. In terms of average monitoring time, C³Net is active for the whole duration of the sequences whereas the other approaches lose the targets in some cases. Finally, the average distance of the target center of mass from the FOV center is reported. By keeping targets closer to the center it is less probable that they will exit the FoV. C³Net on average keeps targets close to the center with the other approaches having a slightly worse performance in this respect due to some bad localization of bounding boxes. In some cases some false positive detections also contributed to the decreased performance as they introduced jitter to the camera motion. Overall, C³Net does not face such challenges as it is not affected by issues related to bad localization of bounding box, errors of the tracker due to retaining outdated bounding box information, multiple overlapping bounding boxes that can cause the camera to remain in an area with less targets, and false positives. The performance gains between C³Net and other methods can be attributed end-to-end nature of the approach that associates lower-level-features with control actions. The continuous frames experiments also allowed to understand how the network behaves under different motion patterns and situations. As shown, in Fig. 8-a when a new target enters the FoV the solution will focus on the center of all visible targets until one of them moves away in which case it will focus on the most prominent target(s) in its view.

4.4 Use-Case Evaluation

To evaluate the performance of the proposed approach in the real-world under different circumstances an embedded platform equipped to a smart camera and a UAV to provide on-board processing. In the experimental setup we use the network output that is related to the camera FoV angles and thus directly controls the camera motion. The same exact

Table 2: Monitoring results on 3 different sequences.

Method	Avg. Number of Targets in FoV			Average Monitoring Time			Avg. of Target Centroid		
	Seq. 1	Seq. 2	Seq. 3	Seq. 1	Seq. 2	Seq. 3	Seq. 1	Seq. 2	Seq. 3
C ³ Net (Proposed)	8	3.3	13	100%	100%	100%	21	18	24
SVM-HOG w tracking	4.4	2.9	8	77%	97%	96%	41.2	33	43.9
tinyYOLO w tracking	5.8	3	10	90%	100%	100%	26.5	26.9	39.06
PDN w tracking	5	3	10.5	98%	74%	100%	25.8	24.2	28.4



Figure 9: Images from the experimental setup evaluation as well as the embedded smart camera (Best viewed in color)

values to move the camera in the simulated environment are used to move the camera in the real world. It is worth noting that none of data captured from the experimental setups were used during the development and training of the network. As such, these experiments also test the generalization capabilities of the network on novel images taken from completely different sources. The real-world experiments also allow for evaluating the processing and timing aspects which was not possible in the simulated environment.

4.4.1 Smart Camera Experiment

C³Net is first evaluated on an embedded smart camera based on the Raspberry Pi computer and a webcam [3]. The webcam is mounted on a motorized two degrees-of-freedom (DoF) pan-tilt stage. The two angular positions are controlled independently using a corresponding servo motor controlled by the Raspberry Pi. Using off-the shelf motor components the time delay to send the camera control and perform the action was very small since the platform interfaces directly with the motor electronics, thus for the purposes of our experiments it is considered negligible [3]. The camera is positioned at a certain height and is tilted with pan range $[-80, \dots, 80]$ and tilt range $[0, \dots, 45]$. The network sends the motion vector to the camera controller that converts it to angles using Eq. 1 and moves the camera head. The indoor environment poses another challenge since there are more reflections, and various background objects in the area. Some results from the lab experiments, as well as the smart camera are shown in Fig. 9. Even in environments which are significantly different from the training set, C³Net manages to follow the target(s).

4.4.2 UAV experiment

To test the algorithm under different conditions we also perform experiments by controlling the camera of a UAV. We assume the same processing platform as with the previous use-case. A series of controlled experiments were carried out with 5 different targets forming a dataset suitable for UAV target detection and following. The purpose of

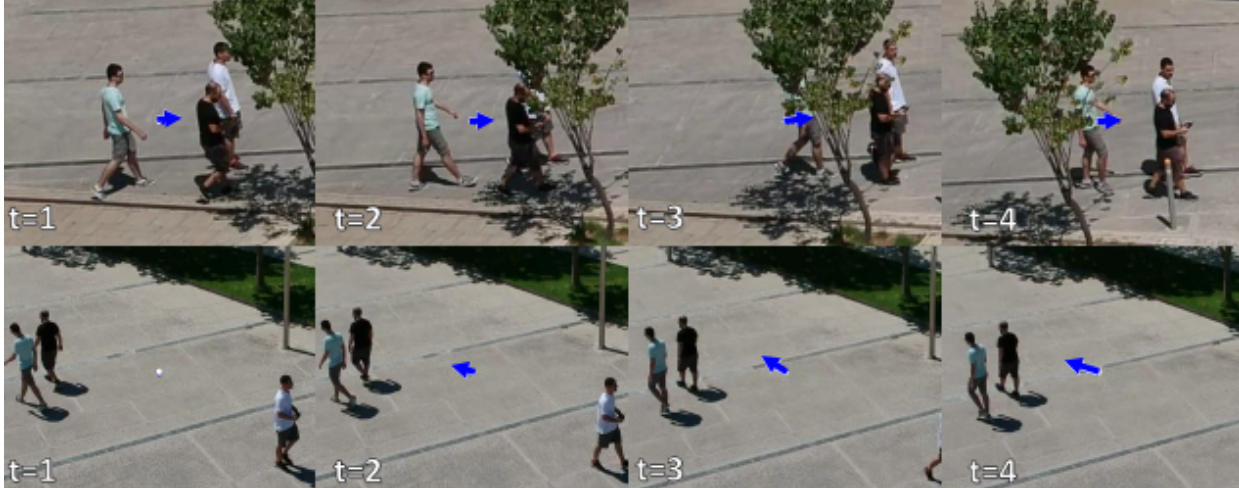


Figure 10: Images taken from the UAV experiment and not used in training. (a)Seq1: The camera using the proposed approach manages to follow the targets even when occluded. (b)Seq2: For a group of targets moving in opposite directions the camera stays with the majority of targets. (Best viewed in color)

the experiments was to observe the behaviour of the algorithm as targets perform different motion patterns, such as moving in opposite directions or move behind an obstacle. Examples from these experiments are shown in Fig. 10. In the first case, the targets move in a group behind a tree. The algorithm is still able to follow the targets even though they are partially occluded. In the second case, the algorithm is centered at the field of view of the targets and as they diverge it chooses to stay with the majority of them since more features are present in that area of the image. Notice that the targets appear in the two sequences at different resolutions but the algorithm is still able to follow them. This also indicates that the algorithm is rather robust to variations in the target-camera distance. In our experiments good performance was observed when the target(s) were in the range of $10 - 50m$, i.e., not too close and not too far for the camera. This suites a variety of smart camera applications.

4.4.3 Benchmarking the computing/processing time

To benchmark the processing time we measure single frame performance which is more representative for real-time camera applications and calculate the resulting FPS. Hence, batching techniques are not applicable. The same processing platform is considered for both use-cases. As shown in Table 3, the proposed approach running on the smart camera platform can achieve a performance of over 10 frames-per-second (FPS) which make it suitable for surveillance applications using embedded smart cameras [42]. The other approaches require considerably more processing time on such low-power platforms (4 FPS at best) and hence are not suited for embedded active smart camera applications. The increased latency also hinders the performance of other approaches since targets may drastically change their direction of movement within the time necessary to calculate a new control output. Further C^3 Net requires $\sim 40\times$ less memory ($\sim 4MB$) compared to tinyYOLO ($\sim 177MB$). The nets are also benchmarked on the Titan-XP GPU platform for reference. Notice, that C^3 Net is an order of magnitude faster than the other approaches.

Table 3: Computing results based on single frame performance

	C^3 Net	PDN	tinyYOLO	SVM-HOG ¹
FPS(embedded)	13	4	3	4
FPS(TitanXP)	367	85	52	N/A
Memory (MB)	4	18	177	N/A

¹No GPU implementation

5 Discussion

Overall, the results are encouraging in that C^3 Net learns the camera control task and generalizes to unseen environments. It first builds an internal representation of the object(s) of interest and learns how to move based on the aggregated feature maps. As such, it has learned to detect stationary targets within a single frame and is guided by their positioning and center of mass.

The proposed end-to-end approach is more robust in finding the control actions as it eliminates the need to rely on bounding box predictions which necessitate complex networks, can be noisy and face difficulties with dense targets as illustrated by the results. Such an approach is well suited for scenarios where no one target is more important than the others but persistent monitoring is necessary while following the maximum number of targets within the FoV. As the objective function is implicitly defined through the training set the approach is flexible and can be adapted to different purposes. For example, the action can be weighted by the distance of the targets from the image center as to give more emphasis on targets at the center of the image that have a higher probability of remaining within the image in future time steps.

Another aspect that is important in practical applications is the response time of the system. The slower it is, the higher the probability that a target will exit its FoV if they suddenly change their motion patterns. The proposed system can tolerate such changes in the order of 0.1 seconds, and can be further improved.

The proposed approach can operate alongside other vision subsystems, can provide robust camera control while other subsystems can perform counting, detection and reidentification at different time intervals. C³Net can be used to follow a group of people while blobs from the activity map, which is computed at no additional cost, can provide regions of interest.

6 Conclusion

This work investigated how to design a more efficient smart camera controller for active vision applications via end-to-end learning using deep convolutional neural networks. A small CNN architecture referred to as C³Net is proposed that maps input images to pan and tilt motion commands. This alleviates the need to have multiple sub-components leading to improved processing times. Even with single image information the proposed network outperforms other multi-stage approaches in terms of monitoring efficiency. Finally, it provides higher frame-rates (a speedup of $\sim 4\times$) while being lightweight validating the assumption that end-to-end approaches can lead to smaller and simpler systems. Our results indicate that when appropriately formulated, end-to-end approaches can lead to more efficient systems thus providing a promising avenue for future research. As future work more elaborate architectures will be explored to take advantage of temporal information. In addition, even though this paper has primarily tackled pan and tilt camera motion, the present study can provide a foundation for similarly controlling other parameters such as zoom.

Acknowledgment

The author would like to acknowledge the support of NVIDIA Corporation with the donation of the Titan Xp GPU used for this research.

References

- [1] C. Micheloni, B. Rinner, and G. L. Foresti, "Video analysis in pan-tilt-zoom camera networks," *IEEE Signal Processing Magazine*, vol. 27, no. 5, pp. 78–90, Sept 2010.
- [2] F. Angella, L. Reithler, and F. Galesio, "Optimal deployment of cameras for video surveillance systems," in *2007 IEEE Conference on Advanced Video and Signal Based Surveillance*, Sept 2007, pp. 388–392.
- [3] C. Kyrkou, E. G. Christoforou, S. Timotheou, T. Theocharides, C. Panayiotou, and M. Polycarpou, "Optimizing the detection performance of smart camera networks through a probabilistic image-based model," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 5, pp. 1197–1211, May 2018.
- [4] R. Wang, H. Dong, T. X. Han, and L. Mei, "Robust tracking via monocular active vision for an intelligent teaching system," *The Visual Computer*, vol. 32, no. 11, pp. 1379–1394, Nov 2016. [Online]. Available: <https://doi.org/10.1007/s00371-015-1206-8>
- [5] H. Chen, X. Zhao, and M. Tan, "A novel pan-tilt camera control approach for visual tracking," in *Proceeding of the 11th World Congress on Intelligent Control and Automation*, June 2014, pp. 2860–2865.
- [6] M. Kiran, V. Tiwari, L. Nguyen-Meidine, L. Morin, and E. Granger, "On the interaction between deep detectors and siamese trackers in video surveillance," in *16th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. Los Alamitos, CA, USA: IEEE Computer Society, sep 2019, pp. 1–8. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/AVSS.2019.8909864>
- [7] H. R. Patil and K. S. Bhagat, "Detection and tracking of moving objects; a survey," *International Journal of Engineering Research and Applications*, vol. 5, no. 11, pp. 138–142, 2015.

- [8] S. V. e. Christophe Bobda, *Distributed Embedded Smart Cameras: Architectures, Design and Applications*, 1st ed. Springer-Verlag New York, 2014.
- [9] M. Bojarski, D. D. Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, X. Zhang, J. Zhao, and K. Zieba, “End to end learning for self-driving cars,” *CoRR*, vol. abs/1604.07316, 2016. [Online]. Available: <http://arxiv.org/abs/1604.07316>
- [10] Q. Bateux, E. Marchand, J. Leitner, F. Chaumette, and P. Corke, “Training deep neural networks for visual servoing,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 3307–3314.
- [11] X. Miao, X. Zhen, X. Liu, C. Deng, V. Athitsos, and H. Huang, “Direct shape regression networks for end-to-end face alignment,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5040–5049.
- [12] J. Ferryman and A. Shahrokni, “Pets2009: Dataset and challenge,” in *2009 Twelfth IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, Dec 2009, pp. 1–6.
- [13] B. Bhanu, C. V. Ravishankar, A. K. Roy-Chowdhury, H. Aghajan, and D. Terzopoulos, *Distributed Video Sensor Networks*, 1st ed. Springer Publishing Company, Incorporated, 2011.
- [14] Z. Zhao, P. Zheng, S. Xu, and X. Wu, “Object detection with deep learning: A review,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 11, pp. 3212–3232, Nov 2019.
- [15] J. Hosang, R. Benenson, and B. Schiele, “Learning non-maximum suppression,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 4507–4515.
- [16] 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 (ICCV)*, Oct 2017.
- [17] C. Feichtenhofer, A. Pinz, and A. Zisserman, “Detect to track and track to detect,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 3038–3046.
- [18] R. P. Pflugfelder, “Siamese learning visual tracking: A survey,” *CoRR*, vol. abs/1707.00569, 2017. [Online]. Available: <http://arxiv.org/abs/1707.00569>
- [19] C. Neff, M. Mendieta, S. Mohan, M. Baharani, S. Rogers, and H. Tabkhi, “Revamp2t: Real-time edge video analytics for multicamera privacy-aware pedestrian tracking,” *IEEE Internet of Things Journal*, vol. 7, no. 4, pp. 2591–2602, 2020.
- [20] C. Ding, B. Song, A. Morye, J. A. Farrell, and A. K. Roy-Chowdhury, “Collaborative sensing in a distributed ptz camera network,” *IEEE Transactions on Image Processing*, vol. 21, no. 7, pp. 3282–3295, July 2012.
- [21] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, “Simple online and realtime tracking,” in *2016 IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 3464–3468.
- [22] Ser-Nam Lim, A. Elgammal, and L. S. Davis, “Image-based pan-tilt camera control in a multi-camera surveillance environment,” in *2003 International Conference on Multimedia and Expo. ICME '03. Proceedings (Cat. No.03TH8698)*, vol. 1, 2003, pp. I–645.
- [23] A. Biswas, P. Guha, A. Mukerjee, and K. S. Venkatesh, “Intrusion detection and tracking with pan-tilt cameras,” in *2006 IET International Conference on Visual Information Engineering*, Sept 2006, pp. 565–571.
- [24] K. Bernardin, F. van de Camp, and R. Stiefelhagen, “Automatic person detection and tracking using fuzzy controlled active cameras,” in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.
- [25] P. S. Dhillon, “Robust real-time face tracking using an active camera,” in *Computational Intelligence in Security for Information Systems*, Á. Herrero, P. Gastaldo, R. Zunino, and E. Corchado, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 179–186.
- [26] T. Dinh, Q. Yu, and G. Medioni, “Real time tracking using an active pan-tilt-zoom network camera,” in *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Oct 2009, pp. 3786–3793.
- [27] M. A. Haj, A. D. Bagdanov, J. Gonzalez, and F. . Roca, “Reactive object tracking with a single ptz camera,” in *2010 20th International Conference on Pattern Recognition*, Aug 2010, pp. 1690–1693.
- [28] W. Luo, P. Sun, F. Zhong, W. Liu, T. Zhang, and Y. Wang, “End-to-end active object tracking via reinforcement learning,” in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 80, 10–15 Jul 2018, pp. 3286–3295.
- [29] M. Al Haj, C. Fernández, Z. Xiong, I. Huerta, J. González, and X. Roca, *Beyond the Static Camera: Issues and Trends in Active Vision*, T. B. Moeslund, A. Hilton, V. Krüger, and L. Sigal, Eds. London: Springer London, 2011.

- [30] Y. Salih and A. S. Malik, "Depth and geometry from a single 2d image using triangulation," in *2012 IEEE International Conference on Multimedia and Expo Workshops*, July 2012, pp. 511–515.
- [31] N. Bo Bo, F. Deboeverie, P. Veelaert, and W. Philips, "Real-time multi-people tracking by greedy likelihood maximization," in *Proceedings of the 9th International Conference on Distributed Smart Cameras*, ser. ICDS '15. New York, NY, USA: ACM, 2015, pp. 32–37. [Online]. Available: <http://doi.acm.org/10.1145/2789116.2789125>
- [32] F. Chollet *et al.*, "keras," <https://github.com/fchollet/keras>, 2015.
- [33] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, "Tensorflow: A system for large-scale machine learning," in *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation*, ser. OSDI'16. Berkeley, CA, USA: USENIX Association, 2016, pp. 265–283. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3026877.3026899>
- [34] D. Chahyati, M. I. Fanany, and A. M. Arymurthy, "Tracking people by detection using cnn features," *Procedia Computer Science*, vol. 124, pp. 167–172, 2017, 4th Information Systems International Conference 2017, ISICO 2017, 6-8 November 2017, Bali, Indonesia. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1877050917329113>
- [35] J. Redmon and A. Farhadi, "Yolo9000: Better, faster, stronger," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 6517–6525.
- [36] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 1, June 2005, pp. 886–893 vol. 1.
- [37] G. Bradski, "The OpenCV Library," *Dr. Dobb's Journal of Software Tools*, 2000.
- [38] V. Campmany, S. Silva, A. Espinosa, J. Moure, D. Vázquez, and A. López, "Gpu-based pedestrian detection for autonomous driving," *Procedia Computer Science*, vol. 80, pp. 2377 – 2381, 2016, international Conference on Computational Science 2016, ICCS 2016, 6-8 June 2016, San Diego, California, USA.
- [39] F. Al Machot, M. Ali, and A. Haj Mosa, "Real-time raindrop detection based on cellular neural networks for adas," *J Real-Time Image Proc*, vol. 16, pp. 1–1, 2019.
- [40] H. Fan and H. Ling, "Siamese cascaded region proposal networks for real-time visual tracking," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 7944–7953.
- [41] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 618–626.
- [42] V. Kulathumani, S. Parupati, A. Ross, and R. Jillela, *Collaborative Face Recognition Using a Network of Embedded Cameras*. London: Springer London, 2011, pp. 373–387.