

Research Result

Real time Multiple Objects Detection using Convolution Network

Kanika Mangla¹, Dr. Vimal Kumar²

^{1,2}Dept. of CSE, Lingayas Vidyapeeth, Faridabad, (U.P.) INDIA

ABSTRACT

Now a day, lots of researchers are working on real-time multiple object detection and tracking algorithms, as the frequent not only in security and surveillance but also in widely applied in various fields, such as health-care monitoring, autonomous driving, anomaly detection, and so on. Multiple Object Tracking is the process of locating multiple objects over time in a video stream. Multiple object Detection is the process of associating detected objects in consecutive video frames. The detected objects may belong to various categories such as vehicles, humans, swaying trees or other moving objects. In this research paper we have proposed an approach to design, implement and evaluate a model that will detect multiple objects in an image, localize the objects and classify them according to their classes.

KEYWORDS

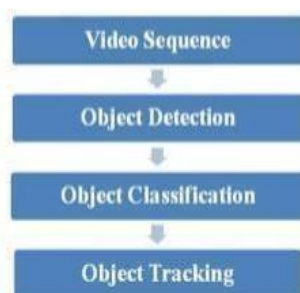
Convolutional neural networks, digital image processing, Multiple object Detection and tracking, Video analysis, Artificial Intelligence

1. INTRODUCTION

Object Tracking

Object detection and classification are two prior steps before performing tracking over video scene. First finding the foreground entities (using features) which are considered as object hypothesis, and then verifying these candidates (using a classifier). Object detection is the process of locating an object of interest in a single frame. So, in other words we can say that object classification is the process to classify these objects using classifier.

The basic flow diagram for multiple object tracking is as following.



Basic Step of Multiple Object Tracking

In this work, we contribute methodology for recognizing objects that does not dispose of any casing in live video, driving consequently the control choices to be made on the most redesigned visual data of the world. Our calculation, Wave3D, accomplishes this component by having the capacity to handle each new edge as a continuation of the preparing of past edges. The calculation produces 3D theories of the nearness of the articles in the genuine

3D world with respect to the robot. These speculations are approved by anticipating them to the 2D picture arranges in the pictures got by the robot camera. The intriguing perspective is that these 3D speculations are free from a particular picture outline, instead of the 2D projections that rely on upon the robot's engines position. Due to this autonomy, our calculation makes utilization of each new picture brought by the camera, which give usto perform the item discovery a chance to handle in live video. Besides, we accomplish another change by producing 3D speculations on the most likely positions in light of the past casings The algorithm advances with its pursuit in a wave-like way, asa spread from the purpose of the speculations. The Wave3 D algorithm extraordinarily enhances the preparing exertion which, liberates calculation for assignment related practices. Visual reconnaissance frameworks have been being used to screen security touchy areas.

2. RELATED WORK

The segmentation of moving objects is an important problem in image sequence analysis and in the problem of video retrieval (i.e. [2, 5, 13]). There is significant related research [1-18]. Some approaches toward visual matching apply stochastic parsing [14]; integrate learning [1, 5, 8]; and adaptive background mixtures [7]. A real-time system for tracking people, called Pfinder ("Person finder"), is proposed by Wren, et al. [16]. [3], a method for spatio-temporal segmentation of long image sequences of scenes that include multiple independently moving objects is presented. This method is based on the Minimum Description Length (MDL) principle.

In literature [5], feature-based algorithm using Kalman filter motion to track multiple objects is proposed by

authors. They had used background subtraction method to detect and extract moving object. Algorithm is validated on human as well as vehicle image sequence and also under confusing situation; it achieves efficient tracking of objects.

Literature [6] focuses on tracking players in football match. Authors have used modified Hungarian algorithm and Kalman filter. As the result concludes that for tracking multiple objects in football match, the linking process is achieved by modified Hungarian algorithm and motion model building and prediction is achieved by Kalman filter successfully.

In literature [7], parallel Kalman filter is used for moving object detection and tracking. The precision and recall value proves that, proposed method is effective for detecting and tracking multiple objects.

In literature [8], authors have proposed unique method named as Dual Layer Particle Filtering (DLPF), which simultaneously detect and track multiple target objects. It uses parent-particles (PP) in first layer to detect multiple objects and child-particles (CP) in second layer to track that detected object.

In literature [9], the proposed method is based on particle filter, which overcomes some challenges of object tracking. It can track the target object robustly when target object is occluded by the background object or other object. Authors had written that, particle filter has higher flexibility than Kalman filter, also they modify particle filter to overcome different challenges.

In literature [10], the proposed method is based on multiple hypothesis tracking. The method is for generic object tracking, which means that there are no priori restrictions in type of objects that can be tracked. However, the authors confess that, the proposed method is quite complex too.

In literature [11], adaptive template matching algorithm is used for tracking the human upper body. The proposed method is fast and robust, as they had added only head edge detection for better tracking. But also, they assume that a person's upper body and face are visible without any occlusion.

In literature [12], authors had written that, in some complex situations like target object have scale changes or similar color with the background, traditional mean-shift algorithm cannot obtain accurate results. So, they suggest new mean shift target tracking algorithm, named as DEPTH & SIFT-Meanshift algorithm. This algorithm is proposed by using a depth camera and SIFT (Scale Invariant Feature Transform) feature metric. The experimental result shows that, the proposed method has abilities to overcome described challenges.

Literature [13] proposes a robust approach for tracking arbitrary objects. Authors propose a new motion model based on Kernelized Harmonic Means and particle filter, and they introduce their proposed model within a SVM framework.

The power of CNNs in image classification was shown already in 1989 when LeCun et al. classified handwritten

digits for zip code recognition with 5% test error [LBD+89]. Dahl et al. show that using ReLUs and dropout improve performance of deep neural networks [DSH13]. Dropout is a method that reduces overfitting by randomly dropping units and their connections during training, thus preventing correlation between units [SHK+14]. In 2015, Ioffe et al. introduce batch normalization, which is another technique of making neural networks more robust[IS15].

Data augmentation is a model-agnostic method for learning invariances. Paulin et al. state the importance of learning class invariances and proposed an explorative algorithm to find the best combination of data transformations [PRH+14].

Paulin et al. show that there are, however, also transformations that do not increase or can possibly even decrease accuracy while increasing computational load [PRH+14]. Another method of enhancing performance in CNNs is to use recurrent convolutional layers as proposed by Liang et al. [LH15].

3. DATA AND RESULTS

Data: Labeled preparing pictures permit us to perform straight relapse, connecting the elements of picture sets (depicted in segment III) with following calculation blunder. Our preparation information comprises of several picture sets (from consecutive video outlines), where the objective item's area is known and marked in both casings.

Acceptance information comparably contains around a hundred picture sets. These pictures are intended to cover an extensive variety of articles and situations, as the preparation is planned to deliver a hearty Kalman channel for any item or environment without "knowing" any exceptional data about it is possible that one.

On a couple of preparing pictures (img_i ; img_{i+1}), we first introduce all following calculations utilizing the objective item's known area in img_i . At that point we figure every calculation's expectation of the item's area in edge img_{i+1} , and record every calculation's blunder in Euclidean separation. The blunder foreseeing highlights for every calculation itemized toward the end of area III are additionally recorded.

After a mistake has been gathered for every preparation pair, we utilize straight relapse on the blunder anticipating components to locate the minimum squares fit to the following blunder. This direct relapse creates a tenet for registering estimation mistake covariance R_k from pictures img_k and image 1. For straightforwardness, we accept that the following calculation blunders are free, so we just n to perform relapse on every calculation's own fluctuation as opposed to a whole covariance lattice. Since physically naming and preparing picture sets is costly and inclined to mistakes, and relapse on such a variety of measurements Naturally produced and named preparing information. The past casing is on the left. requires an expansive arrangement of information, we have supplemented our preparation information with some naturally preparation repairing sets. Fig. 1 demonstrates a case preparing pair, containing an objective item (red circle) on a dark foundation with shaded clamor. spokenk (spoke to as a blue oval) is produced for both pictures. The mark on the main picture is

utilized to instate the following calculations. The name on the second pi re is contrasted with the following calculations' yield (case appeared in white). Our preparation set incorporates around eight thousand such naturally created sets, and our acceptance set incorporates around two thousand.

Training Metrics: Preparing produces a straight capacity for figuring an estimation mistake covariance framework; we apply this standard in a Kalman channel over the preparation and acceptance picture streams to compute the subsequent Kalman channel's following exactness. Preparing and acceptance precision is computed as the rate at which the Kalman channel's anticipated area is inside 48 pixels (Euclidean separation) of the genuine area (gave by name). Precision on both preparing and approval information has been genuinely steady at 97.45%.

Testing

Testing had two parts: Kalman channel vigor in specific situations, and robot sending. All testing is contrasted with the CamShift calculation alone as a benchmark.

1) Kalman Filter Robustness: To test the Kalman channel's power, we ran the Kalman channel on video streams from three hand-picked "troublesome" situations. Our mistake metric for this test is more abnormal state and work escalated than the one utilized as a part of preparing: every video is a trial, and every trial is fruitful just if a human judge chooses that the Kalman channel has carried on accurately.

A. Lighting Changes: These recordings begin by following an article in one lighting environment. The article then moves into an in an unexpected way lit environment, and after that back to the first environment. For instance, the item might be under brilliant surrounding lighting in the principal casing, and after that move against an unforgiving backdrop illumination before at long last coming back to its beginning stage. Achievement is accomplished when the Kalman channel is clearly as yet following the item in the last edge. Our Kalman channel performed inadequately on these test sets, maybe as a result of its substantial dependence on CamShift. It just effectively followed the item 1 in 5 times. A benchmark CamShift tracker comparably succeeded 1 in 5 times.

B. Diverting Background: This set has recordings in which the objective item is comparably hued to the foundation. For instance, the objective might be a dull blue shirt, while the foundation contains a comparably dim blue easy chair. The

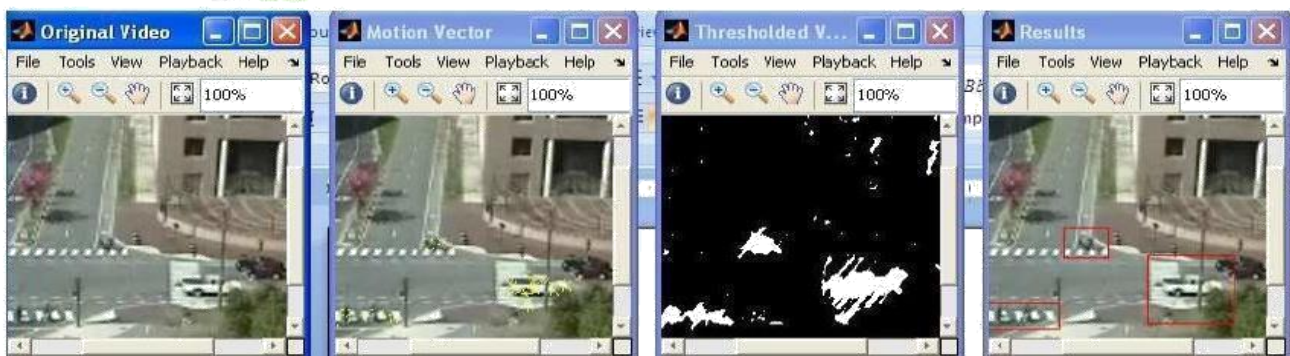
channel is viewed as fruitful on the off chance that despite everything it tracks the article as the item 1 travels through the foundation and back. The prepared Kalman channel was genuinely great on these test sets, succeeding 11 of 13 times. The pattern CamShift alone had more terrible exactness at 8 triumphs out of 13. SURF is especially great at this test set, seeing that it can recognize highlights from comparably hued objects (see Figs. 2 and 3).

C. Clouded Object: These recordings included the objective article moving behind an obstruction and rising up out of the other side, then coming back to its beginning area. Trials are effective if the channel plainly tracks the item as it travels through the hindrance, then back to its beginning stage. Our Kalman channel had a win rate of 8 in 14 on this information set. CamShift performed better at this undertaking, scoring 10 of 14.

D. Robot Deployment: The robot organization tests were genuinely straightforward: every trial comprised of sending the quadrotor running the Kalman channel and controller, and letting it physically track an article.

ALGORITHM

1. import the necessary packages.
2. construct the argument parse and parse the arguments.
3. initialize the list of class labels MobileNet SSD was trained to detect, then generate a set of bounding box colors for each class.
4. load our serialized model from disk.
5. initialize the video stream, allow the camera sensor.
6. loop over the frames from the video stream.
7. grab the frame from the threaded video stream and resize it to have a maximum width of 400 pixels
8. grab the frame dimensions and convert it to a blob. (dnn)
9. extract the confidence (i.e., probability) associated with the prediction
10. extract the index of the class label from the `detections`, then compute the (x, y)-coordinates of the bounding box for the object
11. draw the prediction on the frame and show the output frame.



Multiple object tracking using optical flow from random frame

4. METHODOLOGY & RESULT ANALYSIS

Surveillance means closely and clearly observation of behavior and activities of the objects. Detection of a moving object is necessary for any surveillance system. A static camera can detect and track an object as long as the object is inside the frame of the camera. But as the object goes beyond the boundary of the camera frame, the camera stops tracking it, which is a major limiting factor for the use of a static camera. This limitation can be overcome by using a rotating camera, which will keep continuously rotating and track objects.

CNN detectors and classifiers achieve state-of-the-art performance on a 200-class detection and a 1,000-class classification task with a dataset of approximately 1.2 million training images [RDS+15]. Therefore, we create a dataset pool designed specifically for the task of person classification. In the following, details on the data sources and merging are discussed. Then, the training of a CNN on this data is explained and, lastly, localization based on feature maps is examined.

The entire methodology is divided into 5 basic parts.

- Setting up the object detection directory and virtual environment
- Gathering and labeling dataset
- Creating a label map
- Training & testing the dataset
- Predicting output of new dataset using trained data

Dataset

Datasets such as the ILSVRC and MSCOCO datasets, which are commonly used for deep learning, consist of multiple hundred thousand images. We therefore collect annotated person images from 31 different datasets, creating a large-scale dataset specifically for person classification, which is crucial to training a deep CNN with high accuracy [RDS+15]. The dataset, *Person Data*, comprises about 600,000 images. Examples of different types of person images are shown in Figure.



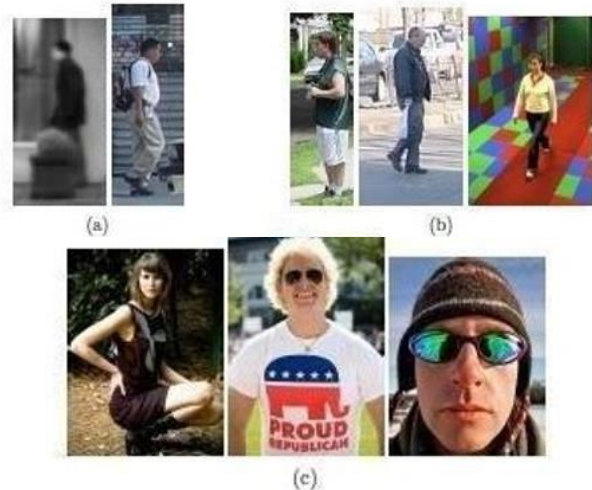
(a) Pedestrian images



(b) General person images

The folder structures, image and annotation formats can be vastly different, resulting in difficult combined usage. To address this issue, we propose a method of extracting

data in a standardized manner from different sources by implementing individual dataset parsers that return standardized information. This requires almost no modification of source datasets and allows for easy integration of new sources as well as simple addition of new queries. The only alterations to source datasets are to decompress.



Examples of image differences among datasets. (a) shows a monochrome and an RGB, (b) different margins and (c) different scales and crops. them, split videos into separate image frames and convert binary MATLAB files to plain text.

The individual parsers handle different folder structures, file naming conventions, and annotation formats, such as XML-based or text file annotations annotation-to-image file correspondence, since one annotation file can include information for only one image or even an entire dataset.



Images extracted from a video of a person walking exhibit little pose variation and, therefore, lead to redundancies.

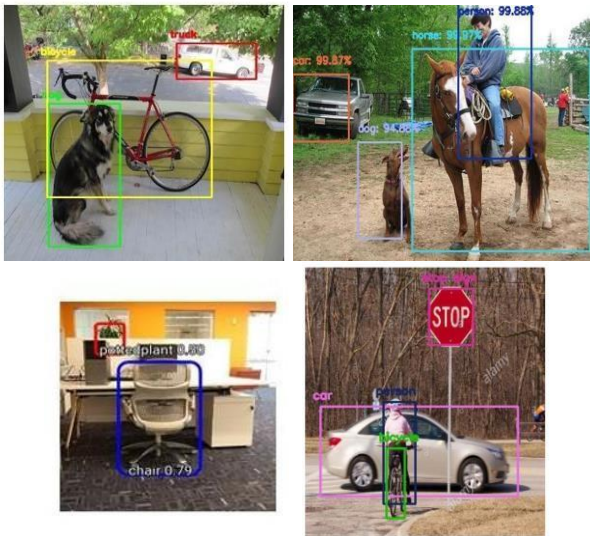
Another script uses the returned standardized information for a specific task, which in our case is image classification. Image classification requires image-level labeling, which means that images containing multiple annotated persons have to be cropped to each annotated bounding box. Therefore, bounding box coordinates and class labels are returned from the parsers and are used to crop and store labeled images in another folder holding the complete, parsed output dataset. Note that bounding boxes can be calculated on the fly from segmentation maps or polygonal annotations, which facilitates the use of segmentation datasets for classification. The results of this procedure can be interpreted as a view of the base dataset pool. In the use of tracking datasets for classification, image redundancy has to be considered.

Video data may show a person walking on a street. Frame-by-frame annotation results in a sequence of images of the same person in a walking cycle. Since the person's pose repeats itself and the background stays relatively similar, e.g. A sidewalk, this means that individual images are similar. Figure 4.3 depicts the small changes in pose variation of a person walking.

Person Classification

Person Data is created as explained in the previous section for the main task of person classification with as much data as possible. A deep CNN is trained to classify an image as containing a person or not. The count of persons in the image is ignored and only the existence of a person is of importance under the assumption that input images only.

5. RESULTS



Muti object detection

6. CONCLUSION & FUTURE DIRECTION

In conclusion our contributions are three-fold. First, we created a large-scale dataset for person recognition with over 600,000 images for person classification. Data was collected from over 30 datasets and a method of viewing such heterogeneous data sources was proposed with respect to their annotation structure, format and type to be able to extract data as available and required for a specific application.

Second, we trained a CNN using this data for binary person classification, i.e., classifying the existence of a person in an image as true or negative, with an error rate of less than 3%.

Third, we demonstrated that subclass labeling aids in training the classifier more robustly and avoiding problems of solver choice and initialization by having stronger labeling of the catch-all negative class. We additionally evaluated the benefit of the data as a pre-training training set for fine-tuned applications.

7. FUTURE WORK

As ImageNet is commonly referenced as a pre-training dataset [RHGS15], possible future work includes comparison of pre-training data with ImageNet or with Person Data for tasks containing people, e.g., detection. Since low-level features are equivalent among different

architectures [LV15a], we suggest that the following multistage fine-tuning could be beneficial. First, the model is trained with ImageNet and learns highly robust low-level features simply due to the immense amount of data. Second, the model is then trained with Person Data to learn mid-level features for person recognition, Lastly, the model is fine-tuned for an application, where application-specific high-level features are learned. As labeling noise is apparent in Person Data, instead of exhaustively, manually pruning the data, it is also possible to adapt the model or the training to explicitly handle it similarly to training with noisy web images [VGLBP15]. Semi supervised methods can be used to identify noisy labels and adapt them. Additionally, prior information about data source that have more labeling noise than others can be utilized to identify and adapt noisy labels. Another way of improving classifier performance is, for example, hyper parameter tuning, i.e., finding the best CNN architecture with layer sizes and properties, and finding the best amount of video redundancy from tracking sources. The latter can, e.g., be done by using image similarity measures to select distinct images from a video sequence. Furthermore, additional data augmentation can be exploited during training but also during testing. Further evaluation can be performed by taking the negative or background class into consideration. The necessity and benefit of separate objects, an additional catch-all class, randomly generated data, and object-less data, i.e., streets, walls, floors, etc., can be analyzed. Finally, Person Data can be used outside of the scope of CNNs.

REFERENCES

- [1] M. B. Blaschko and C. H. Lampert. Learning to localize objects with1 structured output regression. In Computer1 Vision-ECCV 2008, pages 2-15. Springer, 2008.
- [2] L. Bourdev and J. Malik. Poselets: Body part detectors trained using 3d human pose annotations. In International Conference on Computer Vision (ICCV), 2009.
- [3] Cai, Q. Wu, T. Corradi, and P. Hall. The cross-depiction problem: Computer vision algorithms for recognizing objects in artwork and in photographs. arXiv preprint arXiv:1505.00110, 2015.
- [4] N. Dalaland B. Triggs. Histograms of oriented gradients for human detection. In Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, volume 1, pages 886-893. IEEE, 2005.
- [5] T. Dean, M. Ruzon, M. Segal, J. Shlens, S. Vijayanarasimhan, J. Yagnik, et al. Fast, accurate detection of 100,000 object classes on a single machine. In Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on, pages 1814-1821. IEEE, 2013.
- [6] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deepconvolutional activation feature for generic visual recognition. arXiv preprint arXiv:1310.1531, 2013.
- [7] J. Dong, Q. Chen, S. Yan, and A. Yuille. Towards unified object detection and semantic segmentation. In Computer Vision-ECCV 2014, pages 299-314. Springer, 2014.
- [8] Erhan, C. Szegedy, A. Toshev, and D. Anguelov. Scalable object detection using deep neural networks. In Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on, pages 2155-2162. IEEE, 2014.
- [9] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. International Journal of

Computer Vision, 111(1):98-136, Jan. 2015.

- [10] P.F. Felzenszwalb, R.B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(9):1627-1645, 2010.
- [11] Yiwei Wang and John F. Doherty, Robert E. Van Dyck "Moving Object Tracking in Video," 1999.
- [12] Dong Kwon Park, Ho Seok Yoon and Chee Sun Won, "fast object tracking in digital video", *IEEE Transactions on Consumer Electronics*, Vol. 46, No. 3, August 2000.
- [13] Shanik Tiwari, Deepa Kumari, Deepika Gupta, Raina," Enhanced Military Security Via Robot Vision Implementation Using Moving Object Detection and Classification Methods ", *IOSR Journal of Engineering (IOSRJEN)*, Vol. 2 Issue 1, Jan.2012, pp.162-165
- [14] Robert Bodor, Bennett Jackson, Nikolaos Papanikolopoulos," Vision-Based Human Tracking and Activity Recognition", 2000.
- [15] Paul Viola, Michael Jones, Daniel Snow," Detecting Pedestrians Using Patterns of Motion and Appearance ", *Proceedings of the International Conference on Computer Vision (ICCV)*, October 13, 2003, Nice, France.
- [16] Alper Yilmaz, Omar Javed, Mubarak Shah," Object Tracking: A Survey " *ACM Comput. Surv.* 38, 4, Article 13 (Dec. 2006), 45 pages. doi:10.1145/1177352.1177355 <http://doi.acm.org/10.1145/1177352.1177355>
- [17] Lan Wu," Multiview Hockey Tracking with Trajectory Smoothing and Camera Selection ", 2005
- [18] Massimo Piccardi, "Background subtraction techniques: a review ",2004 *IEEE International Conference on Systems, Man and Cybernetics* 0-7803-8566-7/04/\$20.00 @ 2004 IEEE
- [19] Arnab Roy, Sanket Shinde and Kyoung-Don Kang," An Approach for Efficient Real Time Moving Object Detection ", 2009
- [20] Sivabalakrishnan. M and Dr. D. Manjula," An Efficient Foreground Detection Algorithm for Visual Surveillance System " *IJCSNS International Journal of Computer Science and Network Security*, VOL.9 No.5, May 2009.