

A SHORT SURVEY ON FULLY-AUTOMATED PEOPLE MOVEMENT AND IDENTITY DETECTION ALGORITHMS

Maciej Szymkowski¹, Karol Przybyszewski¹

Faculty of Computer Science, Białystok University of Technology, Białystok, Poland

Abstract: Nowadays, diversified companies use security systems based on cameras to increase safety of their enterprise. However, when the camera observes multiple people, it is hard for humans to directly observe each of them. In the literature, there are multiple computer vision-based approaches that automatically detect person identity and the way he is moving. Moreover, there are approaches that identify people across multiple cameras (re-identification). It is crucial, especially in the crowded places. By these algorithms we can detect people whose behavior is strange. Diversified approaches can be easily found in the literature and online-available repositories. The work, presented in this paper, can be divided into three main parts: literature review, selected algorithms implementation and results comparison. We have to claim that each solution was implemented in Python programming language with sufficient libraries. This technology was selected due to its efficiency and simplicity. Results of the conducted experiments have shown that it is clearly possible to detect people's movement and observe their identities even in crowded places.

Keywords: Computer Vision, Image processing, Artificial Intelligence, People movement, Identity detection, Person Re-Identification, Python programming Language

1. Introduction

Nowadays, one of the most important trends in Computer Science is connected with Computer Vision. This technique allows computers to see, identify and process images in the same way as human vision does. The most common algorithms include video and image analysis. In both of these areas, the main goal is to detect a particular object or recognize its identity (as in biometrics algorithms).

Object tracking has been an active field of research for many years. It can be used in video analysis. Its main goal is to keep track of an object's position along the frames, allowing it to preserve its identity. There are plenty of tracking algorithms,

but the main process is to predict the selected object's position in the next frame. This operation is usually less complex (in the terms of computational complexity) than object detection performed in each frame. Multiple cameras provide additional complexity to this algorithm. It is connected with the fact that the solution has to match the same person on many video streams at the same time.

Another interesting approach is connected with direct frame per frame analysis. This method can be compared with the "brute force" algorithm used in password analysis. The similarity can be easily observed due to a couple of facts. The first of them is the high computational complexity of the approaches based on this idea (the same is observed in "brute force" method for password analysis). High usage of resources will not be observed when we analyze one frame, although when the operation is repeated multiple times on a huge amount of frames then the complexity is easily observable. Moreover, if we want to get additional information about the object (for example its identity) we need to use supplementary methods. Probably we will need to connect information from a couple of frames rather than using only one of them.

If it comes to selection of the tracking method from described previously, we can observe that neither of them can guarantee satisfactory results in exceptionally short time. For this aim, we should use another solution that is motion detection. In the case of this approach, we can observe high speed up (in comparison to frame per frame analysis and simple object tracking). We will only observe moving parts of the video, the scene as a whole will not be analyzed (only some specific areas with moving objects, the rest of it will be classified as background). It also should be claimed that there is no single perfect strategy that can guarantee us satisfactory results in a selected environment. In most of the cases, we need to conduct a huge amount of experiments by which we select the proper algorithm.

The described methods can allow us to point some specific goals of Computer Vision. The most important of them is that this technique provides tools used not only to observe selected objects but also to process and return additional information on the basis of observations. Well-known example is a car onboard system for road signs detection and recognition.

The main goal of our work was to find out which algorithms are recently used for movement and identity detection and which of them can guarantee the results in satisfactory short time. We implemented and compared the selection of them. Each experiment was performed on diversified videos. Most samples were collected from DGait [1] and CMU Panoptic Dataset [2] databases.

This work is organized as follows: in the first section the authors describe diversified approaches connected with detection of people movement and tracking. In

the second one, some information about the comparison results are presented. Finally conclusions and future work are given.

2. Related work

Recent advancements in the area of deep learning resulted in transferring those methods into various topics of computer vision. With the success of CNNs [3,4], deep features learned from the networks has replaced handcrafted features [5] for representing person images. One of the important topics is person re-identification (ReID). Given an image of a person captured on one camera, the task is to identify this person from the gallery set captured by other multiple cameras. Table 1 [6] clearly shows an increasing number of papers where deep learning methods were used to develop ReID methods. Table presents number of papers presented at three top conferences:

- CVPR - Conference on Computer Vision and Pattern Recognition [7],
- ICCV - International Conference on Computer Vision [8],
- ECCV - European Conference on Computer Vision [9].

Table 1. The number of deep learning papers related to person ReID included by the three top conferences in recent years.

	2014	2015	2016	2017	2018
CVPR	1	2	5	7	25
ICCV	n.a.	2	n.a.	2	n.a.
ECCV	0	n.a.	2	n.a.	18

Due to the growing efficiency of deep learning methods, we decided to review only papers where those methods were used. Table 2 [6] presents growing accuracy (Rank 1) of the state-of-the-art deep learning models on popular person ReID datasets over the recent years:

- **CUHK03.** The dataset is one of the largest ReID datasets which contains 13,164 images of 1360 identities. All identities are taken from six camera views, and each pedestrian is captured by two cameras. This data set provides two settings. One automatically annotated by a detector and the other manually annotated by humans. Among the two settings, the former is closer to practical scenarios [10].
- **Market-1501.** This dataset consists of 32,643 annotated boxes of 1501 persons. Each pedestrian is collected by at least two cameras and at most six cameras from the front of a supermarket. The boxes of pedestrians are captured by the Deformable Part Model (DPM) detector [11].

- **PRID-2011.** The images of this dataset are captured from two non-overlapping surveillance cameras. One camera captures 749 pedestrians and the other camera captures 385 pedestrians. Among these pedestrians, 200 persons recorded in both cameras. All images are cropped into 128 ×48 pixels. Different from other datasets, PRID 2011 is captured in a relatively clean and simple scene and the dataset has consistent illumination changes [12].

Table 2. Accuracy (Rank 1) of the state-of-the-art deep learning models on popular person ReID datasets over the recent years.

	2016	2017	2018
CUHK03	85,4	88,5	94,9
PRID 2011	66,8	83,7	93,0
Market-1501	83,7	84,9	93,6

Video-based person Re-ID is an extension of image-based person Re-ID. Zheng et al. [13] introduce a large-scale dataset to enable the learning of deep features for video-based Re-ID. They first train a CNN to extract image features then aggregate them into a sequence features with average/maximum pooling. Other works [14] adopt Recurrent Neural Networks to summarize image-wise features into a single feature by exploiting temporal relation within a sequence.

2.1 Capsule networks

Convolutional neural networks (CNNs) have had great success in solving problems with object recognition and classification. However, they are not perfect. If at the input of the convolutional network we give an object in an orientation that the network does not know, or in which objects appear in places that the network is not used to, the prediction task will likely fail. CNN learns statistical patterns on images, but not the basic concepts of what makes something actually look like a specific real object (e.g. a face).

In 2017, Geoffrey Hinton (and others), borrowed ideas from neurobiology that suggest that the brain is organized into modules called capsules [15] (CapsNets). These capsules are particularly good at recognizing features such as orientation (position, size, orientation), deformation, speed, albedo, hue, texture, etc. In the context of neural networks, capsules are represented by groups of neurons.

The results presented in Hinton’s work showed that CapsNets had the highest performance in standard datasets such as MNIST [16] (with a test accuracy of

99.75%) and SmallNORB [17] (with a 45% error reduction over the previous best result). However, the applications and performance of these networks on real and more complex data have not been fully verified. A very important benefit of capsule networks is the transition from black-box neural networks to those that represent more specific characteristics that can help us analyze and understand how the neural network works from the inside.

We should also observe that there are also additional approaches regarding neural networks for object recognition. The most interesting of them was presented in [18]. In this work, the Authors used deep convolutional neural network to classify images from ImageNet LSVRC-2010 competition dataset. What is interesting is that the neural network consisted of 60 million parameters and 650000 neurons. However, the results were not as accurate as expected, the Authors obtained 37,5% and 17% error rates in two testing sets (provided by the organizers of the competition). In the work it was also claimed that the same model was used in the dataset from the another LSVRC competition and the results were much better - error rates were equal from 15,3% to 26,2%. We can conclude that in this case also image quality can have a huge influence on the final results.

What is also interesting is that deep convolutional neural networks and artificial intelligence approaches in general can be used to detect some specific objects. This aim can be mainly observed in biometrics and medical images analysis algorithms. In the first case, neural networks are used to detect specific structures by which human identity can be recognized - for example these can be some parts of fingerprint, eg. not often observed minutiae. In the second idea, machine learning or artificial intelligence is used for detection and segmentation of some pathological changes. The most representative solutions are presented in works [19,20,21]. In these cases, neural networks were used to detect pathological changes in ophthalmic images (especially diabetic retinopathy).

2.2 Selected Computer Vision Surveys

The rapid development of image processing using neural networks has also resulted in a large number of scientific articles describing selected issues in this field. It is worth paying attention to the article "Object Detection in 20 Years: A Survey" [22], where the authors describe extensively over 20 years of history of object detection. The article is based on a review of over 400 papers covering the period from 1990 to 2019. The authors clearly show the division into two main eras of object detection: traditional detection methods (until 2012) and methods based on deep machine learn-

ing (after 2012), among which the most popular concepts are related to convolutional neural networks.

A very good and extensive work describing contemporary deep machine learning architectures is the article "A State-of-the-Art Survey on Deep Learning Theory and Architectures" [23]. The authors comprehensively describe the development of the most important concepts in the field of deep machine learning since 2012. The article also includes a list of the most popular frameworks, SDKs and reference data sets used to implement and evaluate tasks related to deep machine learning.

It is also worth paying attention to the work "A Survey of the Recent Architectures of Deep Convolutional Neural Networks" [24] which focuses on the history of the development of deep convolutional neural networks. The research focuses on showing the internal taxonomy of CNN's latest deep architectures. It also attempts to classify the latest innovations in CNN architectures into seven different categories (English terminology): spatial exploitation, depth, multi-path, width, feature map exploitation, channel boosting, and attention.

3. Conducted experiments

At the beginning of this section, we would like to present the main goal of the conducted experiments as well as the way in which selected algorithms were tested. In the next part of this subsection, comparison between selected methods is also presented.

The main goal of the experiments was to check whether it is possible, with currently used methods, to track people and to gain information about their identity (in biometrics terms). In this case we would like to differentiate people from each other. We do not have data regarding their real identities as well as we do not possess any biometrics databases connected with all analyzed samples.

3.1 Testing procedure

In this subsection we would like to describe how selected algorithms were tested and what was the way to calculate accuracy of the selected solutions. Each algorithm was implemented with Python Programming Language. Testing procedure was realized in the manner described below.

1. In the first step, the selected video was loaded with default Python tools to our algorithm.
2. As the next stage, we marked all people currently visible in the scene. We also differentiate them. Each of them was distinguished with different colors.

3. Further we observed whether the number of marked people was correct. Decision correctness was evaluated on the basis of comparison between the returned number and the number of people observed by the human operator.
4. In the next stage, we changed the scene to the next frame from the camera and observed whether each man has his own, correct marker. It means that we were checking whether marker color was preserved. Once again we compared the algorithm decision with the decision of the human operator.
5. Finally, we combined collected data and evaluated the algorithm. Final accuracy of the solution was calculated as in (1).

$$\eta = 1/|X| \cdot \sum_{x \in X} x_{accuracy} \quad (1)$$

where η is a final, generalized accuracy of the selected algorithm, X is a set of all observations (results of the selected algorithm on each data) whilst $x_{accuracy}$ is an accuracy of the observation x .

3.2 Tested methods and obtained results

In this part of the experiments we used three, most popular algorithms available online: YOLO (You Only Look Once) [25], COCO (Common Objects in COntext) [26] and the default tracking algorithm from OpenCV library [27]. Each of them was analyzed and implemented due to its simplicity and low computational complexity. We assumed that all analyzed solutions can be used in real-time monitoring systems. However, we observed broad differences in the obtained results. This implies that selection of the algorithm can have a huge influence on a final decision. The general scheme of the proposed system is presented in Fig. 1.

Our analysis will be started with presentation of the results obtained with implemented solutions. The captured frame with the marked person is presented in Fig. 2.

It is easily observable that the selected solution can mark a person when it is only one visible. Fig. 3, Fig. 4 and Fig. 5. present behavior of implemented algorithms when there are more people in the scene.

On the basis of the presented results, we can claim that it is clearly possible to observe not only the number of people in the scene but also differentiate them. However, sometimes, selected solutions returned results with mistakenly marked people (e.g. two of them were marked as one - in the terms of identity). This problem is clearly observable in Fig. 4 and Fig. 5. In the first of them, two people were marked as one (algorithm decided that their identities were the same). When it comes to Fig.

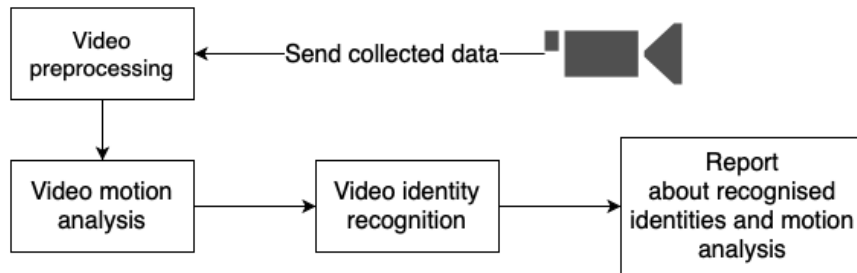


Fig. 1. General scheme of the proposed system.

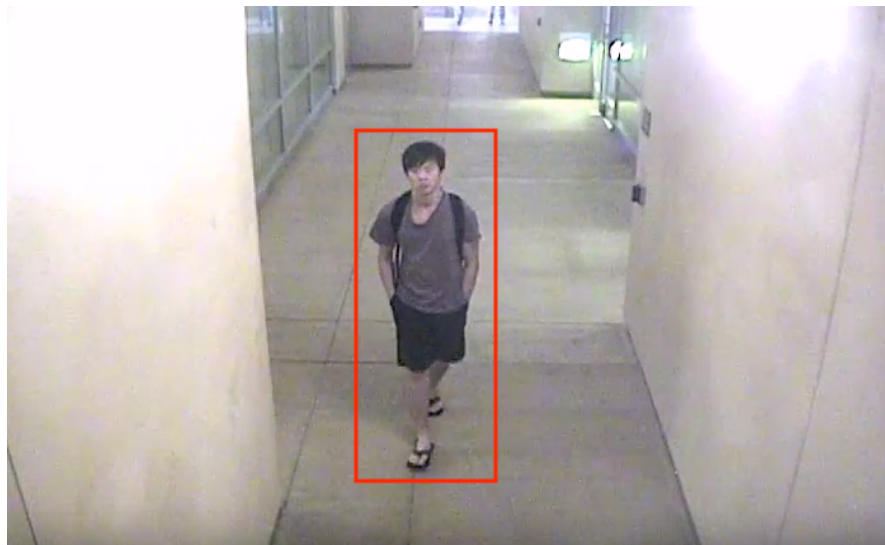


Fig. 2. Image with person marked by implemented algorithms. Frame was taken from database [28].

5 we can observe that two people were not marked. It means that the analyzed solutions do not recognize them as human. These two mistakes were mostly observed. Other types of them mostly appeared only once.

Right now, we would like to present the results of comparison between all tested methods. We made experiments in regards to accuracy (number of people & identities). Of course, we applied slight modifications to have a possibility to obtain information about identities. It is connected with the fact that Standard OpenCV algorithm, YOLO and COCO are the solutions for object detection. Their main goal is to detect objects and to say what is it (man, book, computer... etc.). We applied

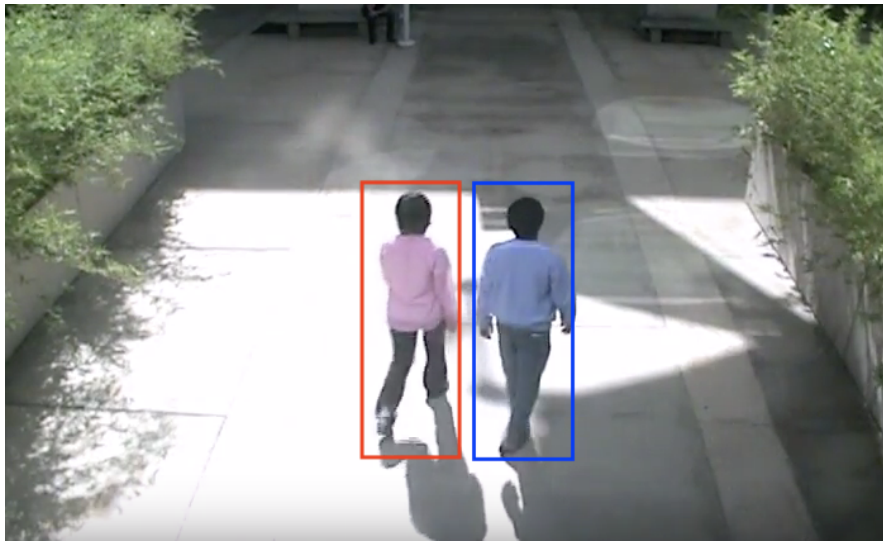


Fig. 3. Multiple people in the scene. Frame was collected from database [28].

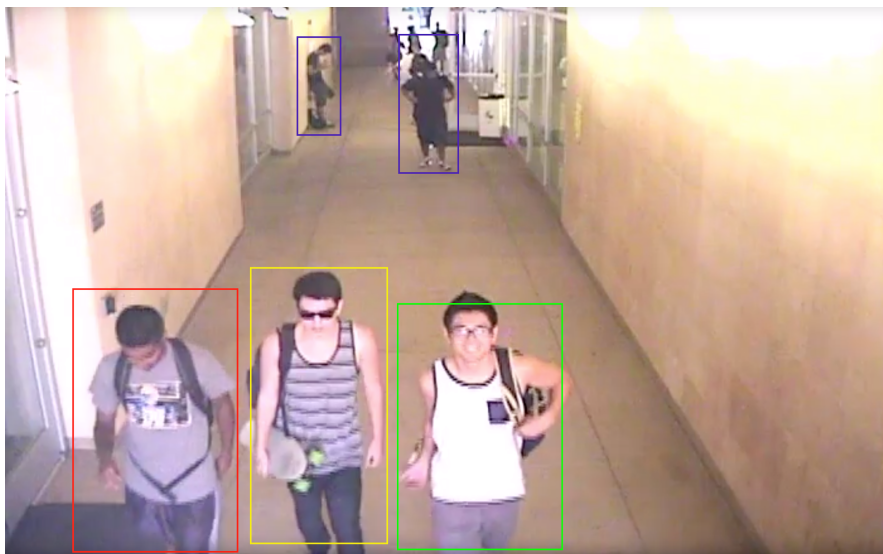


Fig. 4. Multiple people in the scene. Two people were recognized as one (marked with purple rectangle). Frame was taken from database [28].



Fig. 5. Recognition of only two people whilst four appeared in the scene. Frame was captured from the database [28].

some additional changes (regarding information found in diversified guides) for people identification. Each method was tested on more than 100 samples. The results of the tests are presented in Table 3.

Table 3. The results of the experiments.

Algorithm	Mean accuracy (number of people)	Mean accuracy (identities)
Standard algorithm in OpenCV [27]	82,73%	70,2%
YOLO [25]	89,95%	85,36%
COCO [26]	88,92%	81,24%

On the basis of Table 3. It is easily observable that the standard algorithm in OpenCV cannot guarantee satisfactory results even in the case of counting the number of people. YOLO and COCO returned more precise results however neither of them can be used in real circumstances. It is connected with the fact that, in the real environment, the mistake cannot be higher than 0,5%. If we want to replace human

operators by fully-automated algorithms we have to have real confidence that the decision of the algorithm is always right. We cannot risk that it will often fail and the results will be uncertain.

When it comes to identity detection, we observed that each algorithm once again does not return satisfactory results. Once again, the best was YOLO that reached around 85% of correct recognition rate. We observed that the most serious problems with identity detection were generated when the object was partially visible as well as when the light caused its distortion. We think that these problems can be reduced with some additional preprocessing methods by which distortions can be removed and our video frame will be clearer.

4. Conclusions and Future Work

On the basis of the conducted experiments we have to claim that even if it is possible to track people movement and identify them with recently available solutions and tools, the perfect object detection and tracking algorithm has still to be found. The main problem with all analysed approaches is connected with a combination of high speed and high accuracy.

One possible solution is to combine detection and tracking approaches. Analysis of every n -th frame (n has to be selected experimentally) with computationally expensive detection can allow us to update trackers. With this operation we will probably gain a solution that has less computational complexity and can guarantee results similar to the ones obtained with the analyzed approaches. This idea will be implemented and tested in the nearest future.

The Authors' current work is to optimize selected solutions in the terms of time and computational complexity. Moreover, we are working under our own solution for people movement and identity recognition. The experiments in this case are performed on our own database. In the incoming future, we would like to prepare our own fully-automated system based on embedded hardware and FPGAs for people movement and identity detection. We see a huge potential in embedded systems and IoE (Internet of Everything) solutions (with proper cameras and maybe DSP (Digital Signal Processing) modules to obtain precise results in less time.

The performed survey and conducted experiments have shown us that there is still an unresolved gap of problems in video and image processing as well as in object recognition. Moreover, we would like to work under our own re-identification algorithm that will be based on one of the solutions described in this paper. However the selected base has to be perfectly tuned before creation of the final system.

Acknowledgment

This work was partially supported by works WZ/WI-IIT/4/2020 and W/WI-IIT/2/2019 from Białystok University of Technology and funded with resources for research by the Ministry of Science and Higher Education in Poland.

References

- [1] <http://www.cvc.uab.es/DGaitDB/Summary.html> (Access 15.08.2019)
- [2] <http://domedb.perception.cs.cmu.edu/index.html> (Access 15.08.2019)
- [3] Hermans A., Beyer L., and Leibe B.: In defense of the triplet loss for person re-identification, arXiv preprint arXiv:1703.07737, 2017.
- [4] Li W., Zhu X., and Gong S.: Harmonious attention network for person re-identification, In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2285–2294, 2018.
- [5] Cheng D.S., Cristani M., Stoppa M., Bazzani L., and Murino V.: Custom pictorial structures for re-identification, In Bmvc, volume 1, p. 6., 2011.
- [6] Wu D., Zheng S.J., Zhang X.P., et al.: Deep learning-based methods for person re-identification: A comprehensive review, *Neurocomputing* 337, pp. 354–371, <https://doi.org/10.1016/j.neucom.2019.01.079>, 2019.
- [7] <http://cvpr2019.thecvf.com/> (Access 10.11.2019)
- [8] <http://iccv2019.thecvf.com/> (Access 10.11.2019)
- [9] <https://eccv2020.eu/> (Access 10.11.2019)
- [10] Li W., Zhao R., Xiao T., and Wang X.: Deepreid: Deep filter pairing neural network for person re-identification, In Proc. CVPR, 2014.
- [11] Zheng L., Shen L., Tian L., Wang S., Wang J., and Tian Q.: Scalable person re-identification: A benchmark, In Proc. ICCV, 2015.
- [12] Hirzer M., Beleznai C., Roth P.M., Bischof H.: Person re-identification by descriptive and discriminative classification, in: Proceedings of the Scandinavian Conference on Image Analysis, pp. 91–102, 2011.
- [13] Zheng L., et al.: MARS: A Video Benchmark for Large-Scale Person Re-identification, In: European Conference on Computer Vision. Springer, Cham, pp 868-884, 2016.
- [14] Zhou Z., Huang Y., Wang W., Wang L., and Tan T.: See the forest for the trees: Joint spatial and temporal recurrent neural networks for video-based person re-identification, In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 4747–4756, 2017.

- [15] Sabour S., Frosst N. , Hinton G.E.: Dynamic routing between capsules, In Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17), Curran Associates Inc., Red Hook, NY, USA, 3859–3869, 2017.
- [16] LeCun Y., Bottou L., Bengio Y., Haffner P.: Gradient-based learning applied to document recognition, In Proceedings of the IEEE, vol. 86, 2278–2324, 1998.
- [17] LeCun Y., Huang F.J., Bottou L.: Learning Methods for Generic Object Recognition with Invariance to Pose and Lighting, In Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2004
- [18] Krizhevsky A., Sutskever I., Hinton G.E.: ImageNet classification with deep convolutional neural networks, In Commun. ACM 60, vol. 6, pp. 84–90. DOI: <https://doi.org/10.1145/306538>, 2017.
- [19] Wang S., Liang Y., Zhang Y.: Deep Convolutional Neural Networks for Diabetic Retinopathy Detection by Image Classification, In Computers & Electrical Engineering, vol. 72, pp. 274-282, 2018.
- [20] Pratt H., Coenen F., Broadbent D.M., Harding S.P., Zheng Y.: Convolutional Neural Networks for Diabetic Retinopathy, In Proceedings of International Conference on Medical Imaging, Understanding and Analysis 2016, Loughborough, United Kingdom, pp. 1-6, 2016.
- [21] Sarki R., Michalska S., Ahmed K., Wang H., Zhang Y.: Convolutional neural networks for mild diabetic retinopathy detection: an experimental study, DOI: <https://doi.org/10.1101/763136>, bioRxiv, 2019.
- [22] Zou Z., Shi Z., Guo Y., Ye J.: Object detection in 20 years: A survey, arXiv preprint arXiv: 1905.05055, 2019.
- [23] Alom M., Zahangir T.M., Taha M., et. al.: A state-of-the-art survey on deep learning theory and architectures, Electronics vol. 8, no. 3: 292, 2019.
- [24] Khan A., Anabia S., Umme Z., Aqsa Saeed Q.: A survey of the recent architectures of deep convolutional neural networks, Artificial Intelligence Review, pp. 1-62, 2019.
- [25] Redmon J., Farhadi A.: Yolo9000: Better, Faster, Stronger, In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7263-7271, 2017.
- [26] COCO Algorithm, <https://algorithmia.com/algorithms/deeplearning/ObjectDetectionCOCO> (Access 21.08.2019)
- [27] <https://www.pyimagesearch.com/2018/07/30/opencv-object-tracking/> (Access 21.08.2019)
- [28] Zhang S., Staudt E., Faltemier T., Roy-Choudhury A.: A Camera Network Tracking (CamNeT) Dataset and Performance Baseline, In IEEE Winter Conference on Applications of Computer Vision, Waikoloa Beach, Hawaii, 2015.

ANALIZA ALGORYTMÓW SKORELOWANYCH Z DETEKcją RUCHU OSÓB I ICH TOŻSAMOŚCI

Streszczenie Współcześnie w wielu miejscach publicznych oraz obszarach zajmowanych przez zróżnicowane firmy możemy zauważyć systemy bezpieczeństwa bazujące na kamerach. Jednakże bardzo ciężko jest pojedynczemu operatorowi obserwować każdą osobę, która pojawi się na obrazie. W tym celu powstały algorytmy bazujące na metodyce Computer Vision, które mają na celu wykrycie nie tylko trasy poruszania się każdej osoby ale również ocenę jej tożsamości. Co więcej tego typu rozwiązania mogą być bardzo przydatne w zatłoczonych miejscach, gdzie niezwykle ważne jest wykrycie niestandardowego zachowania poszczególnych osób. W literaturze oraz bazach dostępnych online możemy znaleźć zróżnicowane podejścia do rzeczonoego problemu. W ramach naszej pracy porównujemy kilka z nich. Każde z wybranych rozwiązań zostało zaimplementowane przy użyciu języka Python i bibliotek dostępnych w ramach rzeczonoego języka. To środowisko zostało wybrane ze względu na jego wydajność oraz prostotę pisania kodu. Wyniki, które uzyskaliśmy wskazują na to, że aktualnie istniejące solucje mogą być używane do obserwacji trasy poszczególnych osób nawet w zatłoczonych miejscach.

Słowa kluczowe: Przetwarzanie obrazów, Sztuczna inteligencja, Detekcja ruchu, Wykrywanie tożsamości, Język programowania Python

Praca została zrealizowana częściowo na mocy środków pochodzących z pracy WZ/WI-IIT/4/2020 oraz W/WI-IIT/2/2019 przyznanej przez Politechnikę Białostocką z funduszy na badania Ministerstwa Nauki i Szkolnictwa Wyższego w Polsce.