# An Attempt at Automatic Label Generation for Object Entry and Exit on Multimedia with a Semantic Search

Oğuzhan Menemencioğlu [1], İlhami Muharrem Orak [1]

[1]*Department of Computer Engineering, Karabük University, 78050, Karabük, Turkey*

*Abstract –* **Filling the semantic gap between low level keywords which are retrieved automatically from multimedia data and human interpretations of data becomes critical. The research aims to handle the issue by using semantic search on multimedia data. A model is proposed with this research which detects enter/exit points and performs automatic label generation instead of hand using labelling. The model is implemented on a popular and commonly used benchmarking dataset. After the reliability of the model is proofed, it is implemented on a test dataset. It is indicated that the multiple interpretation of results can improve the retrieved information and make the model usability possible.**

*Keywords –* **Semantic search, multimedia, SPARQL, enter/exit point, annotation generation.**

## 1. Introduction

The surveillance systems are widely used in daily life because the development of visual equipment is stunning and the cost-efficient cameras are widespread in streets, shopping malls, airports, etc. Police forces and security agencies consult video records for detecting, identifying and arresting the suspects and the criminals.

On the other hand, the current multimedia data accumulated from different sources are huge and the production rate is very high even when only surveillance systems are considered. Analyzing the current data and retrieving the information is attractive for a large scale of researchers. There is a relation between the multimedia analyzing, required knowledge for an analysis process and the retrieved information after analyzing. The large amounts of domain knowledge is used by high-level activity analysis while very little amounts of domain knowledge is assumed by low-level analysis [1]. The retrieving information process includes three steps: generation of annotation of multimedia, relating them to metadata, and retrieving on these metadata [2]. But the main problem is filling the semantic gap between low level keywords which are retrieved automatically from multimedia data and human interpretations of multimedia data. The semantic gap is defined as disparity between subjectivity, richness of querying, human interpretation and low-level keywords which are extracted automatically from multimedia content. Since the amount of information which multimedia includes (especially in entertainment, security, and teaching or technical documentation fields) is huge, even though the information retrieval of such data sources is very limited, this kind of research becomes very important and valuable [3]. After all, the multimedia annotation generation process can be performed manually or automatically. The most accurate results are obtained by manual annotation generation. In addition, manual annotation generation results are used for the comparison of the automatic generation results as ground truth data [2]. Once the amount of data and the production rate is considered, a critical requirement emerges.

The multimedia content detection or content based approaches are the most popular and important steps for the topic. Thoroughly considered, there are sufficient annotation generation attempts which are ontology based. They use spatial/temporal relations in events and concept definitions to obtain adequate content. However, information retrieval remains in computer vision step. This kind of research only provides successful automatic content extraction

instead of creating content network, and accessing, retrieving and searching the content. They use long time-consuming algorithms such as neural networks, fuzzy, ant colony or genetic for content detection. When the current multimedia data and production rates are considered, a faster approach requirement arises.

Through computer vision research, analysis of multimedia and generation of annotations are already derived. But, when the amount of data is considered, a novel approach is required. Hence, a semantic search model is offered in this context for entry/exit point exploration of objects as a base for furthermore complex annotations on multimedia data.

## 2. Related Works

In the surveillance systems, warning for intrusion or taking timely precaution is provided by many features such as automatic detection of human entry and exit [4]. Automatic detection of human entry and exit can be performed by different ways. There are a few concepts about entry and exit such as points and zones.

An object in surveillance systems may refer to a pedestrian but is not limited to that. Hereinafter, pedestrian or human in multimedia is referred to as object in the paper. An object first appears in entry point and it disappears from the field of view (FOV) in the exit point [5]. Saelao et al. conducted a research on the detection of human entrance and exit. They used the background feature subtraction method (BFSM) and moving foreground feature model (MFFM) which is based on means of FAST features [4]. Video objects have the entry and exit zones which are the locations where objects appear or disappear. The start zones are inferred from the initial position of trajectory and the stop zones are inferred from the final position of trajectory. The zones are modelled by using a mixture of Gaussians (Gaussian mixture model (GMM)) and the zones are learned by using Expectation Maximization (EM) in some of the researches [5]–[10]. Jodoin et al. proposed an approach called orientation distribution functions (ODFs) to recover entry/exit zones [11]. Stauffer proposed a scene-level activity model to define entry/exit zones [12]. These approaches are mostly computer vision based and try to detect entrance/exit in the video processing step. The zones are created by gathering the points of tracking sequences of moving objects.

A Markovian transition model is built by using a probabilistic state space representation in Streib and Davis research. The scene entry/exit locations are estimated by their model, but they did not share comparable results [13]. Park and Trivedi obtained the locations of main entry/exit zones from the frequency of objects' data [14].

An interesting attempt was done by Zhou and Pang. They proposed a holistic approach where the low-level features and metadata are extracted and objects are labeled with their sequential number. The moving information includes position, shape, moving trajectory, duration time, interaction status and queries which are implemented on the metadata [15].

We are inspired to find the entry/exit points from exit point definition of Jodoin et al. in this proposed model. With the proposed approach, instead of using common GMM and EM algorithms, enter and exit points will be obtained by the novel method. By using semantic search, concerned points and more, attractive attributes are retrieved. The details are provided in the following sections.

In the computer vision step, any approach can be preferred for processing the multimedia data. In this research, the background subtraction method is used to process the multimedia. After the video is processed by the favored method, the produced data is questioned by suggested semantic queries on the semantic web infrastructure in an unsupervised form. In other words, a pre-process for zone learning is offered. Details are discussed in the conclusion section.

## 3. Model

For each object, minimum, maximum, and total frame numbers $\sigma_{min_i}$, $\sigma_{max_i}$, $\sigma_{Total\ Frame\ Number_i}$ are calculated with semantic queries on ontologies. For checking consistency, they are interpreted all together and then the objects are labeled.

For consistency, in each video, total frame count is calculated for each object as

$$\sum Total\ Frame\ Number_i = \sigma_{max_i} - \sigma_{min_i} + 1. \quad (1)$$

Frame $\sigma_{min_i}$ can be labeled as "appear" and frame $\sigma_{max_i}$ can be labeled as "disappear". If the result of Eq. (1) and the result of queried $\sigma_{Total\ Frame\ Number_i}$ is same, the frames between $\sigma_{min_i}$ and $\sigma_{max_i}$ can be labeled as "visible". In this case, object $i$ is classified as a true object. Otherwise, object $i$ is classified as false object.

The frame number is detected with the first part of the model. Defining the entry/exit point, the frame number is not qualified alone. To strengthen the definition of entry/exit point, $\sigma_{xc_i}$ and $\sigma_{yc_i}$ are queried with semantic queries on ontologies for each point. Once $\sigma_{t_i}$ is calculated as duration time by semantic queries, the average speed of each object is computed by

$$\sqrt{(\sigma_{yc_2} - \sigma_{yc_1})^2 + (\sigma_{xc_2} - \sigma_{xc_1})^2}/(\sigma_{t_2} - \sigma_{t_1}) \quad (2)$$

using the definition of entry/exit points. Advanced knowledge inferring is attempted by interpretation of entry/exit point definition and the derived speed of the objects.

## 4. Implementation

In this section, the implementation of our model is provided. Details of implementation such as datasets, queries, and interpretations are discussed. Figure 1 represents the overview of the implementation of the proposed model.
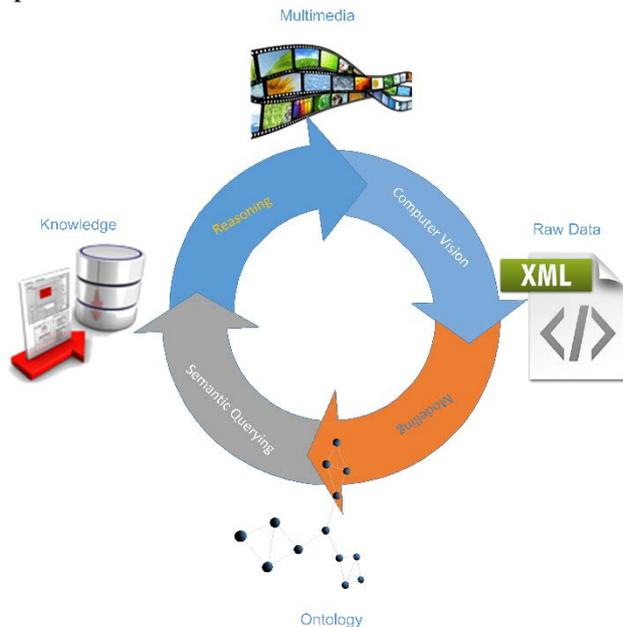


*Figure 1. Lifecycle*

### 4.1. Environment

The mechanism is performed on Microsoft Windows 10, 64-bit running on a system with Intel Core i5-3317U CPU (1.70 GHz) and 4 GB memory. The system can be defined as a personal computer or Ultrabook.

### 4.2. Dataset

In this research, two types of data, benchmarking and testing, are used. Benchmarking data is a public dataset. So, it is ready to use, and in the lifecycle, benchmarking data stays out of the computer vision step. It is modelled to transform it to ontology and get the other steps involved. However, testing data is required to be produced by a computer vision technique.

Benchmarking data is used to compare and assess the model. After the consistency of the model is verified, testing data is used to present the availability of the model.

#### 4.2.1. Benchmarking Data

CAVIAR [16] data set is used for testing the proposed model. CAVIAR is one of the widely-used multimedia datasets. It provides multi-target tracking data with numerous possible entry and exit points, including less crowded scenes and partially occluded people [17], [18].

CAVIAR provides benchmark videos and ground truth files which are labeled by a JAVA-based and hand using interactive tool. Each frame and each object is identified and labelled for comparison with automatic calculations of challengers. In a certain sense, two samples of labeling are provided to explain the meaning. In Figure 2, labeling of the first appearance of the three objects, S0, S1, and S2 are shown in the frame. The labeling of the leaving frames of objects S2 and S3 are presented in Figure 3.

#### 4.2.2. Testing Data

In this study data is obtained from the security cameras of a shopping mall. To generate data for semantic querying from the recordings, any computer vision method could be preferred. Instead of using complex annotation generation approaches, a simple computer vision approach is preferred for this attempt, because it has fast process rate and does not need training etc. Background subtraction algorithm based on the Gaussian mixture processing is used for video analyzing [19].

### 4.3. Queries

The basic semantic queries have been developed in previous paper [20]. The following complicated semantic queries inspired from basic queries suggested in the above-mentioned model are developed thereby to retrieve the targeted knowledge. Each of them could be implemented individually on any benchmarking or test data. Then the results are interpreted together to obtain the target information.

A duration time of FOV denotes the total time that an object appears in between entry and exit points. By running the Query 1, the appearance frame number and the duration time are obtained for each object. Query 2 finds the number of frame and xc, yc values where an object appears first. Query3 gives the result of the number of frames and xc, yc values where an object disappears. In Query 2 and Query 3, the number of objects is limited for presentation of the results because of huge numbers of objects in test data. However, in the benchmarking dataset, the number of objects is around four.

In the queries, the URL of http://semed.karabuk.edu.tr/ns refers to local definition of the data link.

```
# Query1:

PREFIX d: <http://semed.karabuk.edu.tr/ns/data#>
PREFIX                                      fo:
<http://semed.karabuk.edu.tr/ns/frameobject#>
SELECT ?object  (COUNT( DISTINCT?frame) AS
?num_frame) (?num_frame / 25 as ?time)
WHERE {
    ?s fo:object ?object ;
        fo:frameInclude ?frame .
} GROUP BY (?object)
ORDER BY DESC(?num_frame)
LIMIT 10


# Query2:

PREFIX d: <http://semed.karabuk.edu.tr/ns/data#>
PREFIX                                      fo:
<http://semed.karabuk.edu.tr/ns/frameobject#>
SELECT DISTINCT?object ?minFrame ?xc ?yc
WHERE {
    {
        SELECT  ?object  (MIN (DISTINCT?f) as
?minFrame)
        WHERE {
            {
                SELECT   ?object   (COUNT(
DISTINCT?frame) AS ?num_frame)
                WHERE {
                    ?s fo:object ?object ;
                    fo:frameInclude ?frame .
                } GROUP BY (?object)
                ORDER BY DESC(?num_frame)
                LIMIT 10
            }
            ?s fo:object ?object ;
                fo:frameInclude ?frame .
            ?frame fo:frame ?f .
        } GROUP BY (?object)
        ORDER BY (?minFrame)
        LIMIT 10
    }
    {
    ?f fo:frame ?minFrame .
    ?d fo:frameInclude ?f .
    ?d fo:object ?object ;
        fo:xc ?xc ;
        fo:yc ?yc .
    }
} ORDER BY (?object)
LIMIT 10

# Query3:

PREFIX d: <http://semed.karabuk.edu.tr/ns/data#>
PREFIX                                      fo:
<http://semed.karabuk.edu.tr/ns/frameobject#>
SELECT DISTINCT?object ?maxFrame ?xc ?yc
WHERE {
    {
        SELECT  ?object  (MAX (DISTINCT?f) as
?maxFrame)
        WHERE {
            {
                SELECT   ?object   (COUNT(
DISTINCT?frame) AS ?num_frame)
                WHERE {
                    ?s fo:object ?object ;
                    fo:frameInclude ?frame .
                } GROUP BY (?object)
                ORDER BY DESC(?num_frame)
                LIMIT 10
            }
            ?s fo:object ?object ;
                fo:frameInclude ?frame .
            ?frame fo:frame ?f .
        } GROUP BY (?object)
        ORDER BY (?maxFrame)
        LIMIT 10
```

```
    }
    {
    ?f fo:frame ?maxFrame .
    ?d fo:frameInclude ?f .
    ?d fo:object ?object ;
        fo:xc ?xc ;
        fo:yc ?yc .
    }
} ORDER BY (?object)
LIMIT 10
```



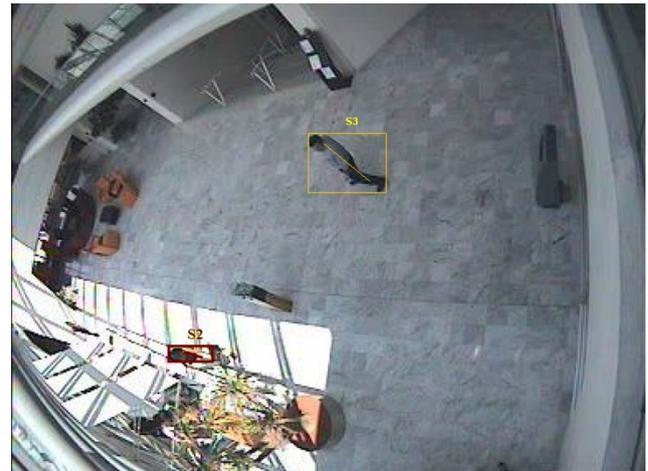*Figure 2. Object appearance in the frame*



*Figure 3. Object disappearance in the frame*

### 4.4. Interpretation

The research focused on semantic web approaches. After the computer vision step, processed data would be available for the proposed model. If the benchmarking data is used, once video processing results are stored in XML, a preprocessing is done for noise removal and the data is modelled in ontologically. If the test data is used, then data is modelled in ontologically. In either case, SPARQL queries are run, and the reasoning process is implemented on ontologies.

For each object entry and exit points are detected. Because of the definition of entry/exit points, the number of frame and xc and yc pixel values for each object are obtained. The total frame number is used

to check the consistency. Furthermore, the visibility duration of each pedestrian could be calculated by the interpretation of the total frame number. The duration time is calculated by dividing the number of frames to 25 frames per second. In interpretation manner, the average speed of objects is calculated by using the values of the xc and yc pixel values and duration time. Then, the xc and yc values of object and the average speed of object is used to infer knowledge by interpretation of all of them together. Figure 4 provides the layer of model.



*Figure 4. Inferring knowledge*

## 5. Results

At first, each query given in Section 4.3 is run on Browse 1 video in Benchmark dataset. The Browse 1 size is 11.82 MB and its length is 42 second. It includes 1042 total frames and 4 objects. For each query, the execution time is included after their results to scale the computation time of proposed queries. The computer used for testing is not fast and powerful enough, and more effective results could be achieved with a computer with better specifications. The specification of the computer used are described in Section 4.1.

### 5.1. The Results of Benchmarking Data

The result of Query1 is shown in Table 1. For each object, the total number of frames, the duration time of FOV is provided.

*Table 1. Total frame number results*

| Object | Num_frame | Time |
|---|---|---|
| 0 | 241 | 9.64 |
| 1 | 244 | 9.76 |
| 2 | 1043 | 41.72 |
| 3 | 709 | 28.36 |

⊕ *Execution time: 0.374 sec*

In Table 2, the frame number of appearance and xc and yc values in entry points are listed. The values of xc and yc indicate the appearance of pixel positions of objects.

*Table 2. The appearance results*

| Object | Minframe | Xc | Yc |
|---|---|---|---|
| 0 | 0 | 287 | 277 |
| 1 | 0 | 210 | 242 |
| 2 | 0 | 84 | 180 |
| 3 | 334 | 26 | 164 |

⊕ *Execution time: 0.483 sec*

The frame number of disappearance, xc and yc values in exit points are listed in Table 3. The disappearance position of objects is indicated by xc and yc values.

*Table 3. The disappearance results*

| Object | Maxframe | Xc | Yc |
|---|---|---|---|
| 0 | 240 | 314 | 286 |
| 1 | 243 | 30 | 168 |
| 2 | 1042 | 114 | 217 |
| 3 | 1042 | 203 | 97 |

⊕ *Execution time: 0.547 sec*

When the results of four queries are considered together, entry/exit points are detected by Query 2 and Query 3. First, the frame number is subtracted from the last frame number and added by 1. If the result equals the value of total frame numbers in Query 1, entry/exit points of the object are accurately detected and in these points (from beginning to end) the object can be labeled as required (appear, visible, and disappear). For instance, an object 0 (0 ID number in Browse1) has been found firstly in frame 0 and frame 240. Per model, total frame number is calculated as 241 by Eq. (1). On the other hand, the total number of frames is obtained as 241 by Query 1. While comparing, the calculated and the obtained total frame numbers for the example, the consistency is derivable. Consistency of the results for benchmarking data is provided in Table 4.

*Table 4. Consistency of benchmarking data*

| Object ID | * | * | ** |
|---|---|---|---|
| | Appear | Disappear | Frame |
| 0 | 0 | 240 | 241 |
| 1 | 0 | 243 | 244 |
| 2 | 0 | 1042 | 1043 |
| 3 | 334 | 1042 | 709 |

* *Number of Frame for labeled with*
** *Total Number of*

Automatic content labeling is achieved by the proposed model. In the implementation, the proposed model is justified with ground truth. For a selected video, moving objects and their actions are required to be identified. In this research, three actions of the objects which are "appear", "visible", and "disappear" are considered. Simply, each object can be determined as "appear" or "disappear" in each frame when it is shown or has left the frame, otherwise it is labeled "visible". It is a simple kind of event detection. The results of labeling on benchmark data are provided in Table 5. The approach is implemented on the twelve ground truth files which are selected among the CAVIAR data set. The amount of the detected objects are provided in the table. The proposed approach achieved almost a hundred percent success. These results show that the proposed approach can be used on any test data. In the data of the video "Walk3", there are five objects. The ID of two and three is the same. After the object two is "disappear" in frame 772, the object three is "appear". So, it is considered that it could not be annotated by experts because it has very small size in the scene up to frame 1097 due to the angle of the camera. It caused a space between frame 772 to frame 1097. If this noise is removed, the accuracy result of the model could be considered at 100 percent success.

*Table 5. Accuracy results*

| Video Name | Count & Percentage | | |
|---|---|---|---|
| | True | False | Pct. (%) |
| Browse_WhileWaiting1 | 3 | | 100 |
| Browse_WhileWaiting2 | 1 | | 100 |
| Browse1 | 4 | | 100 |
| Browse2 | 4 | | 100 |
| Browse3 | 5 | | 100 |
| Browse4 | 3 | | 100 |
| Rest_FallOnFloor | 3 | | 100 |
| Rest_SlumpOnFloor | 4 | | 100 |
| Rest_WiggleOnFloor | 3 | | 100 |
| Walk1 | 4 | | 100 |
| Walk2 | 6 | | 100 |
| Walk3 | 4 | 1 | 80 |
| Average | | | 97.7 |

The average speed is computed considering xc, yc and duration time. Table 6 provides the results of the speed for each object only four videos of benchmark dataset.

It can be seen on the table that although there are various speed values, the speeds of some objects are very similar. The significance of similarity is mentioned in the inferring section.

*Table 6. Speed of objects*

| Object | Browse1 | Browse2 | Browse3 | Browse4 |
|---|---|---|---|---|
| 0 | 2.95233 | 12.62691 | 28.06939 | 31.63714 |
| 1 | 19.94033 | 0.92870 | 10.96453 | 13.36790 |
| 2 | 1.14176 | 2.26637 | 18.27254 | 22.38658 |
| 3 | 6.67336 | 9.69556 | 17.51666 | |
| 4 | | | 18.63390 | |

*Note: Objects of different videos are not same.*

## 5.2. Results of Test Data

After the reliability of the model is proofed on CAVIAR data, Table 7 and Table 8 are presented for test data, which is obtained from shopping mall surveillance camera records. It consists of one hour of records and is provided as an example of test data results.

The consistency of the test data is provided in Table 7. The total number of frames obtained through a SPARQL query proves that the "appear" and "disappear" identifications with other SPARQL queries are correctly done.

*Table 7. Consistency of test results*

| Object ID | * Appear | * Disappear | ** Frame |
|---|---|---|---|
| 2122 | 38816 | 39456 | 641 |
| 2156 | 39500 | 39955 | 456 |
| 2334 | 43047 | 43498 | 452 |
| 2397 | 43988 | 44508 | 521 |
| 2583 | 47404 | 47908 | 505 |
| 2588 | 47436 | 47939 | 504 |
| 3457 | 64687 | 65182 | 496 |
| 3704 | 69945 | 70578 | 634 |
| 3914 | 74338 | 75032 | 695 |
| 4594 | 85510 | 85986 | 477 |

*\* Number of Frame for labeled with*
*\*\* Total Number of*

In Table 8, the frame number of appearance and xc and yc values in entry points; the frame number of disappearance and xc and yc values in exit points are obtained through the queries as shown. Based on this information, the average speed of each object is calculated and listed.

## 5.3. Inferring

When the characteristic of an object's motion is considered, some motions are very similar. In benchmarking data, object 2 and object 3 in video Browse3, have a very similar trajectory. Such information can be obtained by analyzing the results.

When the speed results of objects are stand-alone analyzed, object 0 in video Browse1 and object 1 in

video Browse2 have very similar speeds. Therefore, object 0 in video Browse1, and object 1 in video Browse 2 expose a similar tendency. However, object 1 in video Browse1, and object 4 in Browse3 have very similar speeds but they present very different characteristics. This result implies that taking only speed into consideration does not achieve high accuracy finding similarity of the trajectories. On the other hand, considering speed and appearance/disappearance points (xc, yc, and speed) together indicates more significant information such as similarity of object 2 and object 3 in video Browse3.

When the test data is considered with the mentioned perspective, there is not significant information. But it can be seen from the results that entry and exit positions of the objects vary.

Average speeds, on the other hand are almost around 26 sec/pixel. One object has 35 sec/pixel as the highest value amongst all objects selected. These values could be converted to meter-based values by considering the place where recordings are done.

In the test data, only in one hour of video there are 4849 objects. Therefore, in this application, only top 10 of the most appeared objects in the frames are selected as sample. Future works can focus on the similarity of different features of video objects to analyze, interpret and retrieve the more valuable knowledge.

*Table 8. Results of test data*

| Object | Num_frame | Time | Minframe | Xc | Yc | Maxframe | Xc | Yc | Speed |
|--------|-----------|------|----------|-----|----|----------|-----|-----|----------|
| 2122 | 641 | 25.64 | 38816 | 552 | 1 | 39456 | 202 | 641 | 28.44976 |
| 2156 | 456 | 18.24 | 39500 | 133 | 1 | 39955 | 591 | 456 | 35.3943 |
| 2334 | 452 | 18.08 | 43047 | 564 | 1 | 43498 | 661 | 452 | 25.51512 |
| 2397 | 521 | 20.84 | 43988 | 156 | 1 | 44508 | 380 | 521 | 27.16863 |
| 2583 | 505 | 20.2 | 47404 | 177 | 1 | 47908 | 197 | 505 | 24.97013 |
| 2588 | 504 | 20.16 | 47436 | 369 | 1 | 47939 | 217 | 504 | 26.06471 |
| 3457 | 496 | 19.84 | 64687 | 477 | 1 | 65182 | 113 | 496 | 30.96912 |
| 3704 | 634 | 25.36 | 69945 | 468 | 1 | 70578 | 624 | 634 | 25.70739 |
| 3914 | 695 | 27.8 | 74338 | 273 | 1 | 75032 | 576 | 695 | 27.23962 |
| 4594 | 477 | 19.08 | 85510 | 201 | 1 | 85986 | 10 | 477 | 26.88107 |

## 6. Conclusions

With the proposed model, enter/exit points of pedestrians are identified and labeled as "appear" and "disappear" in these points. Then they are labeled as visible between enter and exit points.

Thus, when the contribution of the proposed approach is considered, automatic label generation is obtained in a simple manner. In the hypothesis of appearance attributes in CAVIAR, there are four values {appear, disappear, occluded, visible}. The framework attempts to generate three of these values automatically on ontologies by querying and reasoning through the proposed model.

Some implications about the characteristics of motion are deduced by using interpretation of the proposed model. Furthermore, the first goal of the researchers extends the paper on analyzing the trajectories by using semantic search infrastructure. Detected entry/exit points will be extended to entry/exit zones by learning from their similarities.

This approach, which is like a simple event detection, may be improved by adding features such as event detection, face or human detection methods in computer vision step displayed in Figure 1. In the dataset, contents such as browse, idleness, walk, interact, reenter, etc. are more attractive ones which still require tackling. The addition of the above-mentioned features will be considered in the following studies.

### References

[1] Morris, B. T., & Trivedi, M. M. (2008). A survey of vision-based trajectory learning and analysis for surveillance. *IEEE transactions on circuits and systems for video technology*, 18(8), 1114-1127.

[2] O. Menemencioğlu and İ. M. Orak, "A Review on Semantic Text and Multimedia Retrieval and Recent Trends," *Int. J. Multimed. Data Eng. Manag.*, vol. 6, no. 1, pp. 54–74, Jan. 2015.

[3] R. Troncy, H. Benoit, and S. Simon, "Introduction," in *Multimedia Semantics*, R. Troncy, H. Benoit, and S. Simon, Eds. West Sussex: Wiley, 2011, pp. 1–5.

[4] Saelao, W., Rattanapitak, W., & Wangsiripitak, S. (2015, October). Tracking-based human entry/exit detection on various video resolutions (A study on parameter effects). In *Information Technology and Electrical Engineering (ICITEE), 2015 7th International Conference on* (pp. 430-435). IEEE.

[5] Makris, D., & Ellis, T. (2002). Path detection in video surveillance. *Image and Vision Computing*, 20(12), 895-903.

[6] Makris, D., & Ellis, T. (2003, July). Automatic learning of an activity-based semantic scene model. In *Advanced Video and Signal Based Surveillance, 2003. Proceedings. IEEE Conference on* (pp. 183-188). IEEE.

[7] Makris, D., & Ellis, T. (2005). Learning semantic scene models from observing activity in visual surveillance. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 35(3), 397-408.

[8] Morris, B. T., & Trivedi, M. M. (2008, September). Learning and classification of trajectories in dynamic scenes: A general framework for live video analysis. In *Advanced Video and Signal Based Surveillance, 2008. AVSS'08. IEEE Fifth International Conference on* (pp. 154-161). IEEE.

[9] Hosseinzadeh, A., & Safabakhsh, R. (2014, February). Learning vehicle motion patterns based on environment model and vehicle trajectories. In *Intelligent Systems (ICIS), 2014 Iranian Conference on* (pp. 1-5). IEEE.

[10] Chen, K. W., Lai, C. C., Lee, P. J., Chen, C. S., & Hung, Y. P. (2011). Adaptive learning for target tracking and true linking discovering across multiple non-overlapping cameras. *IEEE Transactions on Multimedia*, 13(4), 625-638.

[11] Jodoin, P. M., Benezeth, Y., & Wang, Y. (2013, August). Meta-tracking for video scene understanding. In *Advanced Video and Signal Based Surveillance (AVSS), 2013 10th IEEE International Conference on* (pp. 1-6). IEEE.

[12] Stauffer, C. (2003, June). Estimating tracking sources and sinks. In *Computer Vision and Pattern Recognition Workshop, 2003. CVPRW'03. Conference on* (Vol. 4, pp. 35-35). IEEE.

[13] Streib, K., & Davis, J. W. (2010, August). Extracting Pathlets FromWeak Tracking Data. In *Advanced Video and Signal Based Surveillance (AVSS), 2010 Seventh IEEE International Conference on* (pp. 353-360). IEEE.

[14] Park, S., & Trivedi, M. M. (2007). Multi-person interaction and activity analysis: a synergistic track-and body-level analysis framework. *Machine Vision and Applications*, 18(3-4), 151-166.

[15] Zhou, H., & Pang, G. K. (2010, August). Metadata extraction and organization for intelligent video surveillance system. In *Mechatronics and Automation (ICMA), 2010 International Conference on* (pp. 489-494). IEEE.

[16] R. Fisher, "CAVIAR Test Case Scenarios," 2007. [Online].Available: http://homepages.inf.ed.ac.uk/rbf/CAVIARDATA1/.

[17] Yang, B., & Nevatia, R. (2012, June). Multi-target tracking by online learning of non-linear motion patterns and robust appearance models. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on* (pp. 1918-1925). IEEE.

[18] Heili, A., López-Méndez, A., & Odobez, J. M. (2014). Exploiting long-term connectivity and visual motion in CRF-based multi-person tracking. *IEEE Transactions on Image Processing*, 23(7), 3040-3056.

[19] O. Menemencioğlu and İ. M. Orak.(2016). Preprocesses for interpretation of multimedia with semantic search, *Düzce Univ. J. Adv. Technol. Sci.*, 5(2), 131–143.

[20] Menemencioğlu, O., & Orak, İ. M. (2016, May). Semantic querying on multimedia data. In *Signal Processing and Communication Application Conference (SIU), 2016 24th* (pp. 1069-1072). IEEE.