*Research Article*

# Intelligent Analysis for Georeferenced Video Using Context-Based Random Graphs

**Jiangfan Feng and Hu Song**

*College of Computer Science and Technology, Chongqing University of Posts and Telecommunications, Chongqing 400065, China*

Correspondence should be addressed to Jiangfan Feng; fengjf@cqupt.edu.cn

Video sensor networks are formed by the joining of heterogeneous sensor nodes, which is frequently reported as video of communication functionally bound to geographical locations. Decomposition of georeferenced video stream presents the expression of video from spatial feature set. Although it has been studied extensively, spatial relations underlying the scenario are not well understood, which are important to understand the semantics of georeferenced video and behavior of elements. Here we propose a method of mapping georeferenced video sequences for geographical scenes and use contextual random graphs to investigate semantic knowledge of georeferenced video, leading to correlation analysis of the target motion elements in the georeferenced video stream. We have used the connections of motion elements, both the correlation and continuity, to present a dynamic structure in time series that reveals clues to the event development of the video stream. Furthermore, we have provided a method for the effective integration of semantic and campaign information. Ultimately, the experimental results show that the provided method offers a better description of georeferenced video elements that cannot be achieved with existing schemes. In addition, it offers a new way of thinking for the semantic description of the georeferenced video scenarios.

## 1. Introduction

The notion of wireless multimedia sensor networks (WMSNs) is frequently reported as the convergence between the concepts of wireless sensor networks and distributed smart cameras [1]. As a result, an increasing number of video clips is being collected, which has created complex data-handling challenges [2]. Further, some types of video data are naturally tied to geographical locations. For example, video data from traffic monitoring may not contain much meaning without its associated location information. Therefore, most potential applications of a WMSN require the sensor network paradigm to support location-based multimedia services as well as manipulate large scale data at the same time to provide a high quality of experience (QoE), which raises an important issue. How to investigate an intelligent processing method for georeferenced multimedia? Although the question has been extensively addressed theoretically, the method of mapping video sequences to geographical scenes remains to be described. On the other hand, with the growth of geographic information system (GIS) whose major growth area is

the convergence between GIS and multimedia technology, a new paradigm named video-GIS emerged [3–5]. The major researches facing video-GIS are the coding of georeferenced video and the content and types of services that should be provided by georeferenced video. Further improvement of these processes is contingent on deeper understanding of video, as well as improved understanding of the spatial relationship of geographic space. It is due to the necessity of using video-GIS to visualize the relationship between the video analysis methods and the real geographical scene, resulting in georeferenced multimedia intelligent processing method based on context-based random graphs.

Georeferenced video is fundamental process in video-GIS development. Prior research activities on georeferenced video technologies and applications have been conducted. Most of them make use of video and GPS sensors. In [6, 7], Stefanakis and Peterson and Klamma et al. proposed a unified framework for hypermedia and GIS. Pissinou et al. [8] explored topology and direction under the proposed georeferenced video. The work of Hwang et al. and Joo et al. [9, 10] defined the metadata of georeferenced video, which

support interoperability between GIS and video images. In the field of georeferenced video search, Liu et al. [11] presented a sensor-enhanced video annotation system, which searches video clips for the appearance of particular objects. Ay et al. proposed the use of geographical properties of videos [12], while Wang gave a method of time-spatial images to extract the basic movement information [13]. Although single media have been studied extensively, its semantics in geographic space are poorly understood. How to determine the spatial relationship of video elements is one of the most important operations on georeferenced video. For instance, a moving video element changes its position, shape, size, speed, and attribute values over time. Understanding the changing process and rules of these attributes is of important significance to the geographical description of the video.

Many techniques for video event recognition have been proposed. As the work on model-driven methodology which has become well established and approached maturity, the most common and popular conceptualization of fusion systems is the SVM model [14, 15]. However, such methodology not only cannot solve the problems, such as multi-instance, diversity, and multimodal, but needs a large number of training samples. Most previous studies to date have used data-driven method [16] which has been carefully designed to signal clear and distinct semantic of the videos [17–21]. In our event recognition application, we observe that some events may share common motion patterns. Though involved in pattern discovery, data-driven method also contributes to social network during pattern discovery [22–25]. These works have showed a high accuracy in the differentiating of video and its semantic extraction frame. However, most multimedia applications are unknown and uncertain, which are extremely difficult to meet the requirements of real-time stream processing.

Previous studies have shown that multimedia intelligent processing method is important to the development of video-GIS and have achieved inspiring progress. However, these solution methods have suffered from the classical ensemble average limitation presented by the analysis of low-level characteristics. Therefore, the spatial data gathered are sometimes inconclusive and, in part, contradictory. These algorithms usually build or learn a model of the target object first and then use it for tracking, without adapting the model to account for the changes in the appearance of the object, for example, large variation of pose or facial expression, or the surroundings, for example, lighting variation. Furthermore, it is assumed that all images are acquired with a stationary camera. Such approach, in our view, is prone to performance instability, and thus it needs to be addressed when building a robust visual tracker.

To overcome these problems, we will begin by looking at some valid models, which are suitable for georeferenced video understanding and behavior analysis. In this paper, we propose a new event recognition framework for consumer videos by leveraging a large amount of videos. As we know, graph structure provides a complex, dynamic, and robust framework for assembling complex relationships involved in the objects [26], which is suitable for our goal. Thus, multiple random behaviors are presented in certain movement,

making the graph structure unsuitable for describing the real video scenario. To circumvent this problem, random graph model has been taken into consideration, which can be seen as a rather simplified model of the evolution of certain communication net [27]. In our research, it could simplify the analysis of the interaction between video objects substantially for revealing the new insight into the relationships between objects and its complex interaction. Our analysis focuses on describing spatial relationships bound to objects using random graph grammar in georeferenced video, developing a scientific analysis of behavior and structured methods of georeferenced video understanding.

## 2. Preliminary

Surveillance video data is mostly  non-ortho image data so that it does not match up with the geography scene vector data using the traditional method. To solve this problem, a mapping method of video scene imaging data to geography scene vector data is adopted in the paper, as showed in Figure 1. Firstly, the virtual viewpoint camera is constructed by the camera interior and exterior parameters. Secondly, geography scene virtual imaging can be gained from geography scene vector data using the process of model transformation, viewpoint transformation, and pruning according to the computer graphics rendering process, with the corresponding relationship between an object in virtual imaging and vector object. Thirdly, the image matching technology based on the features that have invariant character for translation, scale and rotation is used to match the geography scene virtual imaging and video image. Finally, the corresponding relationship between video image and vector data is established using that between an object in virtual imaging and vector object, with the purpose of accomplishing the mapping of video scene to objects in geography scene.

In the following part, we will introduce several preliminary key steps.

*2.1. Selection Algorithm of Multicamera Based on Spatial Correlation and Target Priority.* Multicamera surveillance system should not only gain detecting and tracking information of motion element of the single camera, but also make the coherent dynamic scene description using all the observations to some extent. Meanwhile, every motion element could be tracked by cameras simultaneously. How to select cameras for tracking a specific target is particularly important in video sensor networks. Based on the spatial correlation [28] and target priority, the paper proposes a selection algorithm of multicamera with task allocation optimized to achieve the automatic selection according to the target priority at each moment.

The algorithm is based on the assumption that a camera with no task carries out the basic single camera tracking which has lower power consumption, and the high-priority task could be preempted when bending. The selection algorithm of multicamera is shown  as in Algorithm 1.

The set of images $I = \{I_1, I_2, \ldots, I_N\}$ is observed by these $N$ cameras, and $S$ denotes the set of cameras selected.
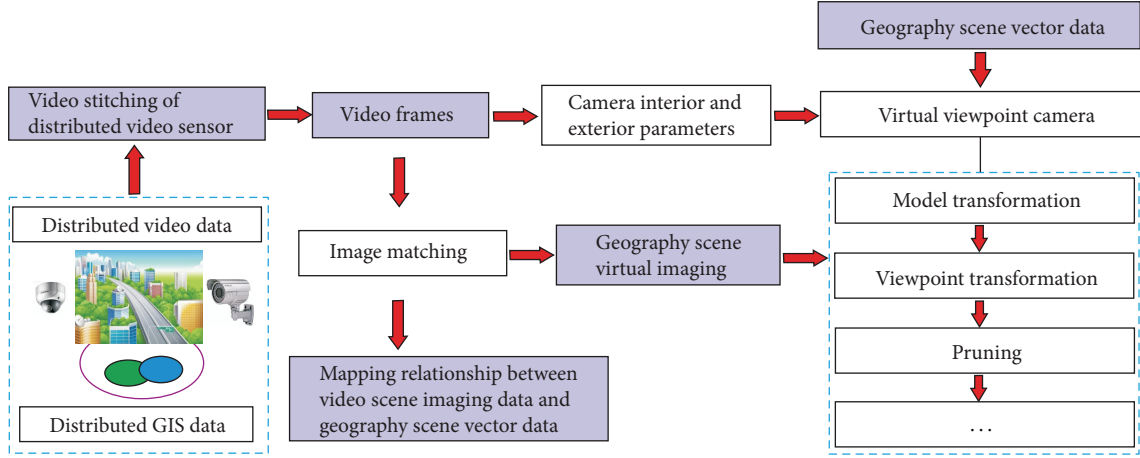
FIGURE 1: Process of the mapping relationship of video scene imaging data to geography scene vector data based on virtual viewpoint.

(1) *begin*
(2) $S = \phi$, $I = \{I_1, I_2, \ldots, I_N\}$, $P_N = P_0$, $\rho(I_i, I_j) = \rho_{ij}$.
(3) *Find* $(I_i, I_j) = \arg \min\limits_{I_i, I_j \in I} \{\rho(I_i, I_j)\}$.
(4) *Add corresponding* $I_i$, $I_j$ *to* $S$. $\{\text{Count} = 2\}$
(5) *for each* $k \in \text{Count}$
(6)       *for each* $(I_{\text{tmp}} \in I,\ I_{\text{tmp}} \notin S)$ *or* $(I_{\text{tmp}} \in I,\ I_{\text{tmp}} \in S, P_{\text{next}} > P_{\text{curr}})$ *do*
(7)             $\rho(I_{\text{tmp}}, S) = \max\limits_{I_p \in S} \{\rho(I_{\text{tmp}}, I_p)\}$
(8)       *end for*
(9)       $I_{\min} = \arg \min\limits_{I_m \in I, I_m \notin S} \{\rho(I_m, S)\}$.
(10)       *add* $I_{\min}$ *to* $S$.
(11) *end for*
(12) *return* $S \subseteq \{I_1, I_2, \ldots, I_{\text{Count}}\}$
(13) *end*

ALGORITHM 1: Selection of multicamera based on spatial correlation and target priority.

$\rho_{ij}$ is correlation coefficient of the two images $I_i$ and $I_j$. The larger the correlation coefficient, the more correlated the two images. In step 6, $P$ denotes the task priority with a default value $P_0$, which can be marked manually by monitoring person. It assigns cameras to the motion element with high priority and coordinates cameras to track different targets based on spatial correlation and target priority.

*2.2. Organization of Video and Location Data.* We have put forward a coding model of video-GIS that is comprised of video and camera's position in conjunction with its view direction and distance. Thus, the location data can be collected automatically by various small sensors to a camera, such as a GPS and a compass (see Figure 2). This eliminates manual work and allows the annotation process to be accurate and scalable. Therefore, we investigate the real-time collection, coding, and integration of video information and GPS information on the SEED-VPM642 platform, and finally we can obtain two different bit-rate location-based streaming media. The lower bit-rate one can be positioned to

the wireless network broadcast live, and the higher one can be positioned to the hard disk storage.

In the coding of video-GIS, we need to calculate the three-dimensional coordinate of the video object [29]. As video-GIS coding based on mobile sensor cannot calculate single video frame by three-dimensional control field, the most effective way is using digital map and spatial geometrical relations (see Figure 3).

Therefore, the geometric relationship among GPS, posture sensor, imaging space, and object space should be built. It is assumed that the axis of imaging space $x, y, z$ is parallel with that of object spatial $X, Y, Z$, respectively. Consider

$$R_G = R_{\text{GPS}}(t) + R_{\text{Att}}(t) \cdot \left[ s_G \cdot R_C^{\text{Att}} \cdot r_g^C(t) + r_{\text{GPS}}^C \right]. \tag{1}$$

In detail, $R_G$ is the coordinate vector of point $G$ in the three-dimensional space. The coordinate function of GPS antenna in the given mapping frame is expressed as $R_{\text{GPS}}(t)$. $R_{\text{Att}}(t)$ represents the rotation matrix function while $s_G$ represents the proportional relationship of image frame and object
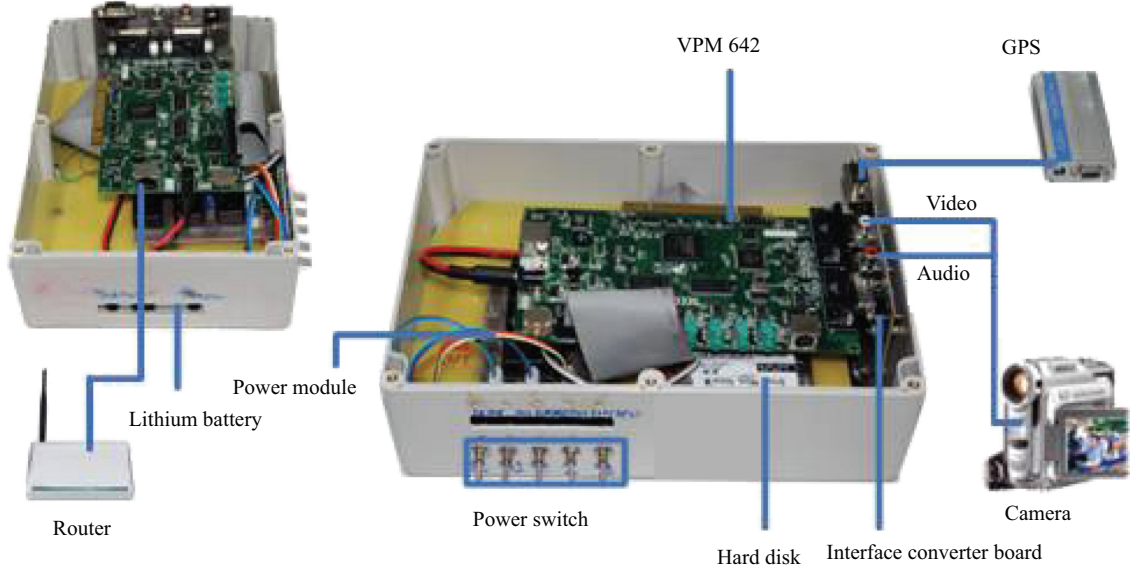
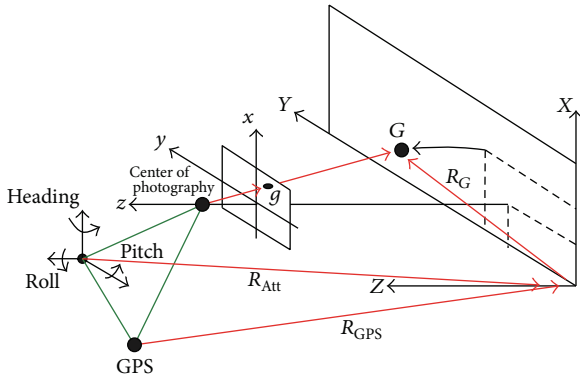FIGURE 2: Experimental hardware and software to acquire georeferenced video.



FIGURE 3: Geometry for calibrating multiple sensors.

TABLE 1: Sample of GPS and MIMU.

| GPS |
| --- |
| UTC 10:12:15 29.564 N 106.585E Alt 213.3 Meters |
| HPR |
| Heading 33.4     Pitch 0.5     Roll 1.3 |

(2) *angle direction elementary information:* including Heading, Pitch, and Roll.

*2.3. Digital Map-Based Image Resolution.* The features of digital maps are expressed by a two-dimensional plane on the vertical projection of the vector data. From the standpoint of this work, the video image is a raster data expression of the feature in the height direction of the information, and video image can also be expressed as the data format of the dotted line surface after the vector processing. Video images and digital map on the point, the line, and the corresponding expression of the surface can be shown at Table 2.

From the view of technology, we subject map-based image resolution to a three-dimensional measurement challenge and then use single-frame video images and digital map matching to define the changes in three steps. The first step is feature extraction of dense range image, which aims to extract the features of point and line. Under the premise of the full calibration to video frame, we can identify the particular characteristics of extracted target to meet the special requirement. For instance, the corners of building or telegraph pole as a fixed line characteristic for the expression of video image is perpendicular to the target. Once formulated, the second step is to combine the line characteristics into the characteristics of the surface using texture information. The third step is matching with digital map vector data. The contents include a variety of different

spatial. Boresight matrix $R_C^{\text{Att}}$ means transformation relation between image frame and main framework of posture sensor. $r_g^C(t)$ represents the vector function of a $g$ point in imaging spatial. And $r_{\text{GPS}}^C$ is the excursion of the geometric center of GPS antenna and the camera lens.

For acquiring a more precise spatial locating information, we need to get the GPS information and attitude information generated by a posture sensor at least. Therefore, the spatial locating information is described by the combination of GPS and angle direction elementary (Heading, Pitch, and Roll), which obtained by Micro Inertial Measurement Unit (MIMU), as shown in Table 1.

As shown in Table 1, there are two kinds of the spatial locating information:

(1) *GPS information:* such as UTC time and longitude latitude;

TABLE 2: Correspondence between video images and digital map.

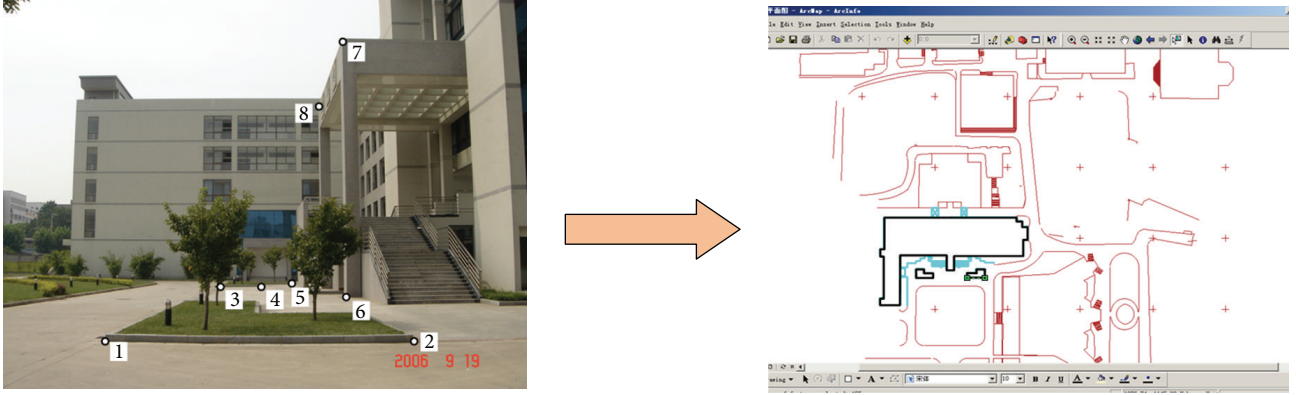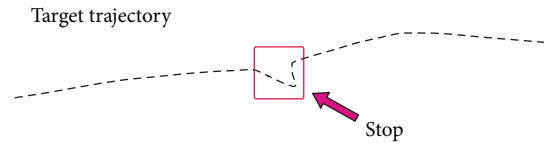| Image → Digital Map | | Digital Map → Image | |
|---|---|---|---|
| Map Symbol | Map Object | Map Symbol | Image Object |
| Point | Point & line | Point | Point & line & Polygon |
| Line | Point & line & Polygon | Line | Point & line & Polygon |
| Polygon | Point & line & Polygon | Polygon | Line & Polygon |



FIGURE 4: Mapping from Image to Digital Map.

matching points, points and lines, a line and a line, and the line and the plane between form and technique, which is shown in Figure 4.

## 3. Syntactic Structure

*3.1. Syntactic Description of Motion Element.* Video motion element mainly refers to the entity objects that could be identified clearly in visual and are important in morphology, such as pedestrians in video surveillance. The description methods of motion element are mainly based on color and texture at present, which is difficult to support the definition of motion element, behavior analysis, and behavior understanding. For a better description of the dynamic characteristic of the video motion element, the paper first gives a definition to some related concepts of motion element.

*Definition 1. State.* *The state* is an abstract of attributes owned by motion element and is a static description of the condition and activity of a motion element at a certain time. *State =* {*Appear, Move, Stop, Disappear*} indicates the basic state of any motion element within the scope of spatial constraint in a georeferenced video stream, including the description information of *Appear, Disappear, Move, and Stop.*

*(a) Appear.* The emerging motion element is newly appear and distinguished from the existing ones in the specific area of geographical boundary, and the state of which is called *Appear.* Then the motion element starts to be detected and tracked. *Appear* instance is regarded as the first instance of motion element.



FIGURE 5: The definition of *Stop.*

*(b) Disappear.* In contrast with the *Appear* state definition, *Disappear* means the state of disappearance in the geographical boundary specific area or the untraceable state within a specific time, which is viewed as the last instance for the state description. *Disappear* state is the signal of canceling motion element detection and tracking.

*(c) Stop.* *Stop* $S$ is defined on triple $S = (Area(S), \zeta_{min}(S), \zeta_{max}(S))$. Among them, $Area(S)$ means the spatial plane area, and $\zeta_{min}(S)$ and $\zeta_{max}(S)$ represent the maximum and minimum time threshold of *Stop,* respectively. And the particular movement or stay that without markedly changed of space coordinate information within a certain region are all viewed as motionless, which is shown in Figure 5.

*(d) Move.* Within the scope of spatial constraint, *Move* $M$ is a general designation of connecting the other three basic states in a continuous motion process of motion element. An instance of *Move* can be represented as $M = (Appear \mid Stop_k, Stop_{k+1} \mid Disappear)$. By connecting the other three basic state instances, *Move* can form a linear sequence formed through the combination of *Appear, Stop,* and *Disappear.*
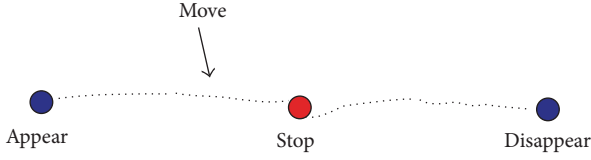
FIGURE 6: *Behavior* state sequence of motion element.



FIGURE 7: A diagram of interaction relation.

*Definition 2. Behavior Attribute.* Behavior description of a single typical motion element mainly includes spatial location and speed. Spatial location can be defined as $Location(Object) = (X_i, Y_i, T_i)$, which means that the spatial location of the motion element *Object* at time point $T_i$ is $(X_i, Y_i)$, and $X_i$ and $Y_i$ represent the horizontal and vertical ordinate values in the two-dimensional plane, respectively. $Speed(Object) = \{S_{Value}, S_{Vector}, T_i\}$ indicates the motion element *Object* with velocity magnitude $S_{Value}$ and velocity direction $S_{Vector}$ at the time point $T_i$, among which $S_{Vector}$ is the unit vector in a general planar domain.

*Definition 3. Relation.* *Relation* is an incidence relation of mutual influence between two motion elements in the same time subspace $T$. $Relation = (Object_i, Object_j, T)$ shows the relationship between motion element $Object_i$ and $Object_j$ in time subspace $T$ which means one-dimensional time coordinates. The measurement of interaction established between the two elements uses probability $P$, which is dynamic adjustment with the influence of temporal-spatial factor, and $P \in [0, 1]$.

*Definition 4. Spatial Relation.* *Spatial Relation* includes measuring relation, direction relation, and topological relation. *Spatial Relation SR* = (*Measure*, *Direction*, *Topology*). *Measure* indicates the measuring relation among motion element using some measure in measuring space, such as distance. In the same planar reference domain, *Direction* is the equity mutual relationship between source target and reference target.

*Definition 5. Visual Feature.* In the georeferenced video stream, the visual characters of one motion element, including *color*, *texture*, and *shape*, will be dynamically changed with the time $T$. Therefore, the changes of visual characters of a motion element within the scope of spatial constraint should be described accurately [30]. And the visual characters mainly include *Color*, *Texture*, *Shape*, and *Size*. *Texture* can reflect the structure mode and gray space distribution formed by local pixels in motion element, while the low-level features can clearly define and describe the motion element.

*3.2. Behavior and Interaction of Motion Element.* In the georeferenced video stream, *Behavior* of the motion element within the specific scope of spatial constraint represents the behavior state sequence, as shown in Figure 6. Let the state set of *Behavior* be a *BehaviorState*, and the typical element is $\tau$ with the definition as follows:

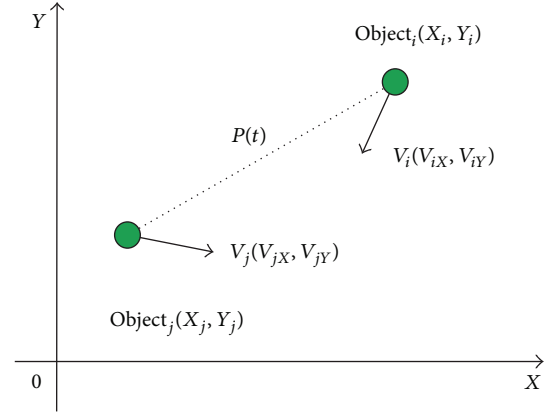$$\tau ::= Appear \mid Move \mid Stop \mid Disappear; \qquad (2)$$

among them, *Appear*, *Disappear*, *Move*, and *Stop* indicate the four basic states of motion element, respectively.

As one of the expression forms of motion element in the video stream, *Interaction* represents the mutual influence or joint action caused during the course of the *Relation* of two behavior state instance. The necessary condition for establishing the interaction relationship is the two incidence relation between the two behavior state instances that exist at the same time. It can be defined as five-meshes

$$Interaction = \{Object, BehaviorState, SR, T, Rule\}. \qquad (3)$$

Under the influence of temporal subspace $T$ and *spatial relation SR*, *Interaction* is the description of mutual influence between motion element $Object_i$ and $Object_j$. Behavior state of *Object* can be any state instance in the *BehaviorState* collection, and interaction production rule and interaction optimization update rule are involved in *Rule*. Therefore, the measuring of interaction has two influence factors, temporal and spatial factors.

Due to the close correlation of spatial relation at any time point $T_{i+1}$ and former $T_i$, the spatial relation at $T_{i+1}$ is always closely related to that at former time point $T_i$. Thus, the spatial relation evolution process among motion elements can be defined as a Markov chain in the temporal subspace $T$, with its evolution having Markov quality

$$P_T \{G_{t+1} \mid G_t, G_{t-1}, \ldots, G_0\} = P_T \{G_{t+1} \mid G_t\}. \qquad (4)$$

Meanwhile, the measuring value $P$ of interaction between the two motion element established *Relation* can be computed based on the planar spatial distance *Distance*, velocity magnitude, and direction angle, including the current topology at time point $T_t$, as shown in Figure 7.

In the georeferenced video stream, the dynamic update function of interaction relation within the scope of spatial constraint is shown as follows:

$$P(t+1)$$
$$= \text{Min}\left[1, \text{Max}\left(0, \sqrt{P^2(t) + \omega(t+1) \times \eta(1 - c(t))}\,\right)\right]. \quad (5)$$

Among them, $P(t)$ represents the interaction relation measuring value between a certain motion element and others, with the range of $P \in [0, 1]$. The higher value indicates the more hospitable relationship. When the interaction is established by behavior state instances, the initial value works as $P(0) = \rho_1 \times Distance(i, j) + \rho_2 \times \theta(i, j)$. $\omega(t)$ indicates the duration of interaction relation with the current state, and the dynamic change of $c$ parameter is shown as follows:

$$c(t+1) = c(t) + a \times \frac{D_{t+1}}{D_t} \times \left(1 - \left|\frac{\omega(t+1)+1}{2} - P(t)\right|\right)$$
$$\times \min\left[c(t), 1 - c(t)\right]; \quad (6)$$

$c(t)$ represents a new confidence level while $\alpha$ learning rate. $\min(c, 1-c)$ ensures parameter $c(t) \in [0, 1]$.

# 4. Semantics and Formalization of Georeferenced Video

For the accurate description and behavior understanding of motion elements in the georeferenced video stream, the paper proposes an analysis method based on sparse random graphs with the purpose of observing the character evolution with time and presents an indicating and measuring method of video motion element with dynamic topology structure information based on context-sensitive sparse random graph grammar.

*4.1. Formalization of Georeferenced Video.* Random graph $G = (V, E, \Omega)$ is defined on triple, while the edge set $E$ of graph $G$ with the vertex set $V$ is defined in probabilistic spaces $\Omega$. Consider

$$P(e_{ij} \in E) = P_{ij}, \quad P_{ij} \in (0, 1), \quad \sum P_{ij} = 1. \quad (7)$$

Each edge of random graph $G$ is mutually independent; namely, any two vertexes that established incidence relation connected independently with probability $P$. As the spatial relation will be dynamically changed during the movement with the time factor, it is necessary to describe the motion state and interaction relationships within specific spatial area using random graph. Context-sensitive sparse random graph grammar can be defined as five-meshes

$$G = (S, V_N, R, \delta, Ch). \quad (8)$$

Among them, $S$ is the root vertex that an initial vertex of semantic event in the georeferenced video stream. There is only one $S$ vertex in the video event sequence. *Vertex* $V_N = \{V_1, V_2, V_3, \ldots\}$ involves all the motion elements emerged in the specific spatial area. $R$ in the formula means the evolution process and rule of random graph $G$ while $\delta$ the state transition functions. The cohesion of random subgraph $Ch$ indicates the inner coupling degree of motion element group.

The motion element vertex of random graph can be defined as follows:

$$V_i = (index, Time, State, Location,$$
$$Speed, Interaction, SR, VF). \quad (9)$$

It shows the motion status and interaction information of a motion element $V_i$ labeled index at the time point *Time*. Among them, *Location* and *Speed* represent the position coordinate and the velocity of motion element $V_i$ in the planar area, respectively. *Interaction* is the description of interaction while $SR = (Measure, Direction, Topology)$ the spatial relation existed in the motion element. *Virtual feature VF* shows low-level features information of a motion element including *Color*, *Shape*, and *Size* at the time point *Time*. *State* $= \{Appear, Move, Stop, Disappear\}$ is the basic state of motion element.

*4.2. Evolution Rule.* As a posterior method, dynamic process of motion elements in the video stream can be visually described and showed based on sparse random graph. The temporal and spatial evolution model of motion element is able to describe the basic character and dynamic process of spatial relation accurately. The essence of dynamic evolution process of sparse random graph is the continuous transition process of state space in random graphs.

Therefore, the state transition function of sparse random graph can be defined as a mapping relation

$$\delta = \Theta \longrightarrow \Theta. \quad (10)$$

Among them, $\Theta$ is the state space of sparse random graph, $\Theta = (d_1, d_2, \ldots, d_n)^T$, and $d$ is a variable in state region.

The dynamic evolution process of sparse random graph includes its character update of motion element vertex $V_N$, emerging vertex with the *Appear* and *Disappear* behavior states, and the dynamic adjustment of edge set $E$ and interaction relation $P$ of random graphs. For the accurate description of event development process in georeferenced video stream, evolution rule algorithm of sparse random graph is shown in Algorithm 2.

We can get the corresponding dynamic evolution model of sparse random graph using the evolution rule algorithm. Step (2) in the algorithm shows the creating and adding root vertex S, and $G_0 = (V_0, E_0) := (\{S\}, \emptyset)$. Adding a new motion element vertex $V_{tmp}$ in sparse random graph $G_{Active}$ is in step (5) while deleting the vanish vertex $V_i$ and its association edge in step (11). Among them, function $getRestriction(V_j)$ in step (18) and $getAttract(V_j)$ in step (20) indicate whether it can delete or add the edge that vertex $V_j$ associated, respectively. Step (27) accomplishes the dynamic update of interaction relation $P$ in sparse random graph $G_{Active}$.

Input: sparse random graph $G_{Active}$, motion element detection and recognition information;
Output: return $G_{Active}$;
(1) IF $t = 0$ Then
(2)　　Create first node $S$ & Add $S$ to $V_N$;
(3) End IF
(4) While $t \geq 1$ do
(5)　　IF $V_{tmp} \rightarrow$ State Is Equal *Appear* Then
(6)　　　　Find nearest node $V_{near}$;
(7)　　　　Create new edge $E(V_{tmp}, V_{near})$;
(8)　　　　Add $V_{tmp}$ to $V_N$;
(9)　　End IF
(10)　For $V_i \in V_N$ do　　//Update all Nodes in $G_{Active}$
(11)　　　IF $V_i \rightarrow$ State Is Equal *Disappear* Then
(12)　　　　　Remove $V_i$ from $V_N$;
(13)　　　　　Delete edge of $V_i$ in $G_{Active}$;
(14)　　　End IF
(15)　　　Update $V_i$;
(16)　End For
(17)　For $V_j \in V_N$ do　　//Update all Edges in $G_{Active}$
(18)　　　IF $Flag \leftarrow getRestriction(V_j)$ Then
(19)　　　　　Delete edge of $V_j$ in $G_{Active}$;
(20)　　　Else IF $Flag \leftarrow getAttract(V_j)$ Then
(21)　　　　　Add new edge of $V_j$ to $G_{Active}$;
(22)　　　End IF
(23)　End For
(24)　For $V_k \in V_N$ do　　//Update $P$ of Graph in $G_{Active}$
(25)　　　IF $V_k \rightarrow$ State Is Equal *Appear* Then
(26)　　　　　$P \leftarrow P(0)$;
(27)　　　Else Update other $P$ of $G_{Active}$;
(28)　　　End IF
(29)　End For
(30) Return $G_{Active}$;

ALGORITHM 2: Evolution rule algorithm of sparse random graph.

### 4.3. Random Subgraph.

*Cohesion* of random subgraph refers to the close relation of motion element. To measuring close relation, the paper introduces the concept of structural entropy. As a measuring method of messiness and randomness of the state, structure entropy is related closely to the compactness of random subgraph. The higher the compactness is, the lower the structure entropy value will be.

If vertexes $V_i$ and $V_j$ have close correlation with each other, then $P(V_i, V_j) = P(t)$. Let $N(i) = \sum_{j=1}^{n} P(V_i, V_j)$, associative strength $\xi(i) = N(i)/\sum_{j=1}^{n} N(j)$. The structure entropy of random subgraph is $H = \sum_{i=1}^{n} \xi(i) \ln \xi(i)$, and $\sum_{i=1}^{n} \xi(i) = 1$. Therefore, the Cohesion of random subgraph is $Ch(G') = -\sum_{j=1}^{n} (N(j)/n) \times (\xi(i) \ln \xi(i)/\ln n)$, with $Ch(G') \subseteq [0, 1]$.

### 4.4. Early Warning of Video Event.

Using the numerical calculation method of interaction relationship, abnormal behavior and emergency in video can be distinguished based on random graph grammar, and the possible special situation can be early warned. There are two different threat levels
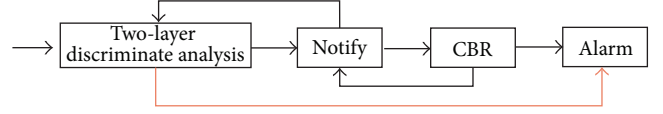


FIGURE 8: Notify and Alarm processing of video event.

generated by video event: notify and alarm, which is shown in Figure 8.

The paper is mainly to detect the unexpected crowd incident and conflict in the massive video events and proposes a novel two-layer discriminate method, which consists of individual attribute layer and group attribute layer. Once occurring video abnormal event, the corresponding real-time status of random graph must be described, which can be expressed as follows.

*(1) Individual Attribute Layer.* The owned velocity of multiple random graph nodes has modified radically in per unit of time $T$, and the relevant movement direction has also changed significantly.

Specifically speaking, the detection and selection of variation range or interval of movement attributes in random graph can use sliding window. In the continuous movement attribute value $V = \{V_1, V_2, \ldots, V_n\}$ in time series, $V_1$ exists before the emergence of $V_2$, while $V_2$ exists before $V_3$. The difference is obtained by the two continuous attribute values. In the paper, the data in the sliding interval $\Delta T$ is viewed as the discriminative and forecasting sample, when the continuous difference $D(V_i, V_j, T)$ is larger than the given threshold, and the sliding intervals $\Delta T$ is within the max time threshold. Otherwise, recalibrate over the entire sliding intervals for new computation.

*(2) Group Attribute Layer.* The multiple interaction and distance values among random graph nodes in groups fluctuate greatly, or the multiple numerical variations of interaction relationship in random subgraph are changed significantly. The discriminant analysis of video abnormal event is achieved according to the check whether the change rate of parameter value $\vec{r}$ is greater than the given threshold $pThreh$, as

$$\vec{r} = \frac{dp}{dt} \geq pThreh. \tag{11}$$

Once either circumstance occurred, it must be entering the next notify phase.

When entering the notify discriminative phase, the random subgraph showing diffusion or flocking status makes numerical calculation. Using the computing method of structure entropy value, the corresponding random subgraph status is measured, and the entropy value $Ch(G')$ is viewed as the warning degree of video abnormal behavior and emergency. With regard to different levels of urgency and security, the warning degree $Warning(t)$ is set to different threshold intervals as follows:

$$Warning(t) = Ch(G') = -\sum_{j=1}^{n} \frac{N(j)}{n} \times \frac{\xi(i) \ln \xi(i)}{\ln n}. \tag{12}$$

The warning degree *Warning*(*t*) is divided into three warning threshold intervals in the paper, which are Warning1, Warning2, and Warning3. Specifically, Warning1 indicates the early warning degree, which means that video abnormal event will be occurred in the next unit time and the discriminative module obtains alertness. Warning2 shows the probable warning degree and is the identifying processing transformed into the CBR phase. If the entropy value of random subgraph is greater than the max value of given threshold interval, the CBR discrimination phase works. Based on the video event features, the traditional CBR method is used to further identification. Warning3 expresses the confirmed warning degree, which can enter the Alarm phase of video abnormal event directly without the traditional CBR method.

The discriminate method based on the random graph is defined as graph-based reasoning (GBR) in the paper, while the improved GBR fused with traditional CBR method is GBR-C. The intelligent analysis for different video scenes plays an important role in the real-time detection of video abnormal behaviors and mass incidents. The instantaneous status information of video motion element is integrated with the random graph model and summarizes the random subgraph patterns and behavior rules with a statistical description. In violation of the behavior regularity of common video events, it is a latent exceptional event, and extracts the features of video motion elements involved which are recorded in object layer stream for the efficient retrieval of content-based video.

## 5. Experiment and Analysis

In order to verify the feasibility and availability of the proposed framework, space information of a motion element is extracted at real-time based on the detection and tracking [31, 32]. According to the dynamic change situation of space semantics, a timing description method using random graph grammar depicts the event development of video stream clearly.

*5.1. Interaction Description. Interaction* is the mutual incidence relation among motion element. For the accurate description of the dynamic change process of interaction relation, interaction *P* should be calculated real-time based on the spatial information in experimental video including planar spatial distance, velocity magnitude, and direction angle. And the calculation results of real-time interaction update function $P(t)$ of the video clip trim from frame 550 to frame 685 is shown in Figure 9.

In Figure 9, function $P_1$ shows a changing trend of increasing first and then decreasing gradually in the video clip. The minimum value of interaction $P_1$ is at frame 685 with the value 0.11 while the maximum is at frame 586 with the value 0.38. And function $P_2$ indicates the changing process of two close targets. The minimum value of $P_2$ is at frame 592 with the value 0.23 while the maximum is at frame 685 with the value 0.79. The increasing planar spatial distance *Distance* and motion direction variation of two motion elements make the

decreasing interaction value. On the contrary, as the planar spatial distance decreases and the duration of interaction continues to increase, interaction value *P* increases gradually.

The previous results show that it can accurately depict the dynamic varying changes of the interaction relation of video motion elements. However, the accurate depiction is an indispensable premise for the description of the georeferenced video stream.

*5.2. Georeferenced Video Stream Description.* Based on the richer spatial semantic of motion elements in the georeferenced video stream, we can realize the intelligent parsing of georeferenced video content using context-sensitive sparse random graph grammar. The spatial relationship of motion elements in image space is transformed to that of object space, and the motion status and interaction relation can be depicted using random graph. The continuous transition process of inner state space in random graph is enforced with the dynamic evolution process of sparse random graph.

With the spatial reference data, the sparse random graph evolution processing based on the monitoring target is achieved. And the consecutive people emerged within the video surveillance range are labeled as A, B, C, and D which are shown in Figure 10. As soon as the moving object appears, a new random graph node will express it; when it leaves the surveillance confine, the corresponding node will disappear while the edge set constituted by the interaction that associated with the node is set to null. Using our video test data, the evolutionary process and timing evolving description diagram of the video clip trim from frame 1041 to frame 1712 is shown in Figure 10.

We can see that the timing evolving description diagram can be constructed by the automatic intelligent analysis and calculation of a video clip, and it verifies the correctness and effectiveness of the evolution rule algorithm of sparse random graph. Within the scope of the specific geographical space, the time-varying attributes of random graph nodes are visual displayed, such as behavior state, spatial location, and movement parameter. And the basis recorded information of each video motion element is shown in Algorithm 3.

Among them, the basic information consists of attribute information, spatial location information, and other movement parameter, which are shown in Algorithm 3. The attribute information *State* indicates the behavior status of the video motion element with succinct expressional number 0, 1, 2 and 3, which are described respectively with the four basic behavior *state* {*Appear*, *Disappear*, *Move*, *Stop*}. And the interaction relationship attribute including the index of two elements, the numerical calculation value of interaction, and the relative spatial directions. The whole structural description of video motion element generated automatically is shown in Figure 11.

The automatically generated file mainly consists of two parts: the configuration data and content data. The movement status information about motion element *Object* in the georeferenced video stream is described in detail in the content data part while the basic attribute information about testing video clip in configure data part. In the continuous period

FIGURE 9: Dynamic change process of interaction relation.



FIGURE 10: Timing evolving description diagram of the georeferenced video stream.

TABLE 3: Three test sample videos.

| Test samples | Alarm types | Time (s) | Scenes number |
|---|---|---|---|
| A | Cross-border | 372 | 18 |
| B | Flocking | 1423 | 42 |
| C | Conflict | 588 | 27 |

of time series, movement status information of each motion element including the behavior state sequence, real-time spatial location information, and the statistical information about interaction relation can be queried directly from the XML file. It also provides a novel simple nonlinear indexing for the understanding and description of video content.

*5.3. Performance of Video Event Warning.* To validate the proposed early warning method of video abnormal behavior and emergency, we analyzed the performance of various attributes using the video test data which involves a crowd video scene. Experimental analysis mainly contains the real-time warning entropy value of random subgraph, warning degree, and real-time changes of corresponding subgraph node number and

```xml
<?xml version="1.0" encoding="UTF-8"?>
<SRG>
    <configData>
        <descriptor name="VideoMotionInfo" type="FILE">
            <attribute Medianame="TrackElement_Session_1"></attribute>
            <attribute FrameCount="76283" FPS="30" ></attribute>
            <attribute CamGPS="GPS" UTC="10:12:15" LAT="29.564" LON="106.585"></attribute>
            <attribute CamMIMU="HRP" Heading="33.4" Pitch="0.5" Roll="1.3"></attribute>
        </descriptor>
        <descriptor name="IMGLOCTInfo" type="LOCAT">
            <attribute ReferPoint="0" PixelX="359" PixelY="251" LocatX="582" LocatY="870" ></attribute>
            <attribute ReferPoint="UL" PixelX="184" PixelY="104" LocatX="0" LocatY="0" ></attribute>
            <attribute ReferPoint="UR" PixelX="346" PixelY="142" LocatX="580" LocatY="291" ></attribute>
            <attribute ReferPoint="DL" PixelX="240" PixelY="337" LocatX="289" LocatY="1164" ></attribute>
            <attribute ReferPoint="DR" PixelX="612" PixelY="321" LocatX="1160" LocatY="1162" ></attribute>
        </descriptor>
    </configData>
    <ContentData>
        <object framespan="1:1500">
            <attribute name="MotionElement">
                ...... ......
                <index="17" State="2" frame="302" timeDelay= 302 PixelX="212" PixelY="171"
                    LoctX="160" LoctY="486" DeltX="0.43" DeltY="0.29" Speed="(0.43,0.29)"
                    InteractionNum="1" Interaction="{(18,17, NE, 0.42)}" VF="0" Other="0"/>
                <index="18" State="2" frame="302" timeDelay= 302 PixelX="374" PixelY="259"
                    LoctX="606" LoctY="898" DeltX="0.33" DeltY="0.72" Speed="(0.33,0.72)"
                    InteractionNum="1" Interaction="{(17,18, SW, 0.42)}" VF="0" Other="0"/>
                ...... ......
            </attribute>
        </object>
        <object framespan="1501:3000">
            <attribute name="MotionElement">
                ...... ......
            </attribute>
        </object>
    </ContentData>
</SRG>
```

Figure 11: Structural description of video motion feature.

the total graph node number, which are shown, respectively, in Figure 12. And the horizontal axis indicates the video running time with 10 seconds as a scale unit.

As can be seen from the previous illustration, the warning entropy value of real-time random subgraph using the computing method of structure entropy value is due to random fluctuations in Figure 12(a). According to the warning degree of video abnormal behavior and emergency, three different warning threshold intervals are set in our test. And the Warning2 degree occurred between 252 and 270 seconds shown in Figure 12(b). The Warning1 indicates the early warning degree in most of the time, which means that video abnormal event will be emerged. Figure 12(c) shows the real-time nodes number of random subgraph in the video surveillance scope while Figure 12(d) shows the total graph node number.
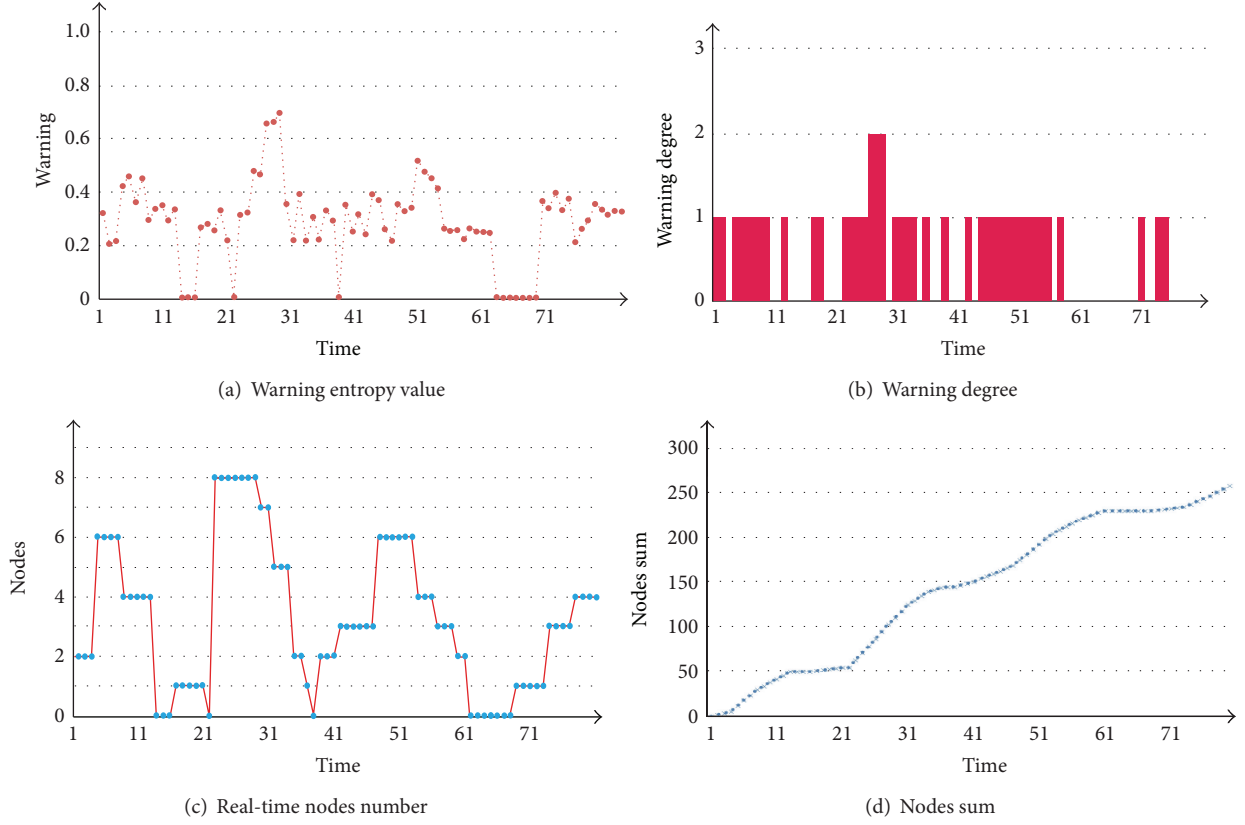
(a) Warning entropy value



(b) Warning degree



(c) Real-time nodes number



(d) Nodes sum

FIGURE 12: Structural description of video motion feature.

```
⟨attribute name="MotionElement"⟩
    ⟨ index="63"                              //Sequence Number
      State="2"                               //Behavior State
      frame="612"                             //Current Frame Number
      timeDelay="612"                         //Duration
      PixelX="198" PixelY="211"               //Image Space Coordinate
      LoctX="45" LoctY="60"                   //Object Space Coordinate
      DeltX="0.85" DeltY="0.24"               //Relative Distance
      Speed="(0.85, 0.24)"                    //Speed
      InteractionNum="1"                      //Interaction Relationship Number
      Interaction="{(64, 63, NE, 0.51)}"      //Interaction Relationship
                                              //(Objecti, Objectj, Direction, P(t))

      VF="0"
      Other="0"/⟩
⟨/attribute⟩
```

ALGORITHM 3: Basic information of each video motion element.

*5.4. Performance Comparisons of Intelligent Analysis Methods.* In this section, we compare the proposed method with other methods, such as the Coarse-Grained SVM, Fine-Grained SVM [15], and MKL [19]. Using the three sample videos (Table 3) which involve some events that contain a group of people interact with each other, we carry out the comparison study. And all the chosen samples are considered as the labeled training data within the target domain.

GBR accomplishes a concise numerical calculation and avoids the problems of computing complexity in the traditional CBR method. In Tables 4, 5, and 6, we compare the performance of GBR, GBR & CBR, with other methods using three different videos.

From Tables 4, 5, and 6 , we observe that GBR extends the processing time in a common detection of video event, but the forecasting accuracy of video abnormal behavior and emergency increased significantly with lower computation

TABLE 4: Comparison of crossing sample A with different methods.

| Method | Out-Detection | Correct-Detection | Omit-Detection | Time (s) |
|---|---|---|---|---|
| GBR | 20 | 18 | 0 | 0.72 |
| GBR (with CBR) | 25 | 18 | 0 | 1.35 |
| Coarse-grained SVM | 40 | 16 | 2 | 1.47 |
| Fine-grained SVM | 26 | 15 | 3 | 1.21 |
| MKL | 30 | 16 | 2 | 1.50 |

TABLE 5: Comparison of flocking sample B with different methods.

| Method | Out-Detection | Correct-Detection | Omit-Detection | Time (s) |
|---|---|---|---|---|
| GBR | 48 | 35 | 7 | 2.74 |
| GBR (with CBR) | 51 | 39 | 3 | 4.33 |
| Coarse-grained SVM | 90 | 37 | 5 | 5.41 |
| Fine-grained SVM | 45 | 35 | 7 | 4.76 |
| MKL | 43 | 36 | 6 | 5.53 |

TABLE 6: Comparison of conflict sample C with different methods.

| Method | Out-Detection | Correct-Detection | Omit-Detection | Time(s) |
|---|---|---|---|---|
| GBR | 35 | 23 | 4 | 3.46 |
| GBR (with CBR) | 36 | 24 | 3 | 4.17 |
| Coarse-grained SVM | 53 | 24 | 3 | 3.90 |
| Fine-grained SVM | 37 | 21 | 6 | 3.54 |
| MKL | 40 | 22 | 5 | 3.95 |

and complexity. Therefore, the energy consumption of sensors will be reduced which is consistent with the transmission costs, especially in the nonrecurring flocking emergency with complex video event modeling.

## 6. Conclusion

In summary, findings from the present study are all based on low-level visual features, which mean that there was a shortage of spatial constraints and coupling analysis with geography environment. It is necessary to establish the relationship between video analysis method and the real geographical scene. A georeferenced video analysis method is proposed based on the context-based random graph. The data are obtained using a wireless network of environmental sensors scattered at the supervising area and a vision sensor monitoring the same geographical area. Experimental results prove that the proposed description method of georeferenced video using random graph is feasible and efficient. Through the intelligent parsing of the georeferenced video data stream, we can get a novel visual description method using random graph which can clearly depict the development clue of video scenes and also offer the possibility to browse the video stream quickly. Meanwhile, random graph can be used as an effective nonlinear indexing for the content-based video indexing and browsing application.

As a future work, we will propose the enhancement of the implemented algorithms with alternative combination rules and the fusion of audio and video to deal with the uncertainty, imprecision, and incompleteness of the underlying information. In addition, large amounts of data should be conducted to set various parameters, such as thresholds, false alarm rates, and fusion weights.

## References

[1] I. F. Akyildiz, T. Melodia, and K. R. Chowdhury, "A survey on wireless multimedia sensor networks," *Computer Networks*, vol. 51, no. 4, pp. 921–960, 2007.

[2] S. A. Ay, R. Zimmermann, and S. H. Kim, "Relevance ranking in georeferenced video search," *Multimedia Systems*, vol. 16, no. 2, pp. 105–125, 2010.

[3] T. Navarrete and J. Blat, "VideoGIS: segmenting and indexing video based on geographic information," in *Proceedings of*

*the 5th AGILE Conference on Geographic Information Science*, pp. 1–9, April 2002.

[4] C. Larouche, C. Laflamme, R. Lévesque, and R. Denis, "Videography in Canada: georeferenced aerial videography in erosion monitoring," *GIM International*, vol. 16, no. 9, pp. 46–49, 2002.

[5] N. Davies, K. Cheverst, K. Mitchell, and A. Efrat, "Using and determining location in a context-sensitive tour guide," *Computer*, vol. 34, no. 8, pp. 35–41, 2001.

[6] E. Stefanakis and M. Peterson, *Geographic Hypermedia: Concepts and Systems*, Lecture Notes in Geoinformation and Cartography, Springer, 2006.

[7] R. Klamma, M. Spaniol, M. Jarke, Y. Cao, M. Jansen, and G. Toubekis, "A hypermedia afghan sites and monuments database," *Geographic Hypermedia*, pp. 189–209, 2006.

[8] N. Pissinou, I. Radev, and K. Makki, "Spatio-temporal modeling in video and multimedia geographic information systems," *Geo-Informatica*, vol. 5, no. 4, pp. 375–409, 2001.

[9] T. H. Hwang, K. H. Choi, I. H. Joo, and J. H. Lee, "MPEG-7 metadata for video-based GIS applications," in *Proceedings of the IEEE International Geoscience and Remote Sensing Symposium (IGARSS '03)*, vol. 6, pp. 3641–3643, July 2003.

[10] I. H. Joo, T. H. Hwang, and K. H. Choi, "Generation of video metadata supporting video-GIS integration," in *Proceedings of the International Conference on Image Processing (ICIP '04)*, vol. 3, pp. 1695–1698, October 2004.

[11] X. Liu, M. Corner, and P. Shenoy, "SEVA: sensor-enhanced video annotation," in *Proceedings of the 13th Annual ACM International Conference on Multimedia*, pp. 618–627, November 2005.

[12] S. A. Ay, R. Zimmermann, and S. H. Kim, "Relevance ranking in georeferenced video search," *Multimedia Systems*, vol. 16, no. 2, pp. 105–125, 2010.

[13] J. Wang and D. Yang, "A traffic parameters extraction method using time-spatial image based on multicameras," *International Journal of Distributed Sensor Networks*, vol. 2013, Article ID 108056, 17 pages, 2013.

[14] C. J. C. Burges, "A tutorial on support vector machines for pattern recognition," *Data Mining and Knowledge Discovery*, vol. 2, no. 2, pp. 121–167, 1998.

[15] Y. Zhuang, Z. Fu, Z. Ye, and F. Wu, "Real-time recognition of explosion scenes based on audio-visual hierarchical model," *Journal of Computer-Aided Design and Computer Graphics*, vol. 16, no. 1, pp. 90–97, 2004.

[16] B. Taskar, P. Abbeel, and D. Koller, "Discriminative probabilistic models for relational data," in *Proceedings of the 18th Conference on Uncertainty in Artificial Intelligence*, pp. 485–492, Morgan Kaufmann, 2002.

[17] C. Wang, L. Zhang, and H. J. Zhang, "Learning to reduce the semantic gap in web image retrieval and annotation," in *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 355–362, July 2008.

[18] R. Fergus, Y. Weiss, and A. Torralba, "Semi-supervised learning in gigantic image collections," in *Proceedings of the Neural Information Processing Systems*, pp. 522–530, 2009.

[19] G. Mehmet and A. Ethem, "Multiple kernel learning algorithms," *Journal of Machine Learning Research*, vol. 12, pp. 2211–2268, 2011.

[20] J. Liu, J. Luo, and M. Shah, "Recognizing realistic actions from videos "in the Wild"," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1996–2003, June 2009.

[21] L. Duan, D. Xu, I. Tsang, and J. Luo, "Visual event recognition in videos by learning from web data," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 9, pp. 1667–1680, 2012.

[22] X. Jin, A. Gallagher, L. Cao, J. Luo, and J. Han, "The wisdom of social multimedia: using Flickr for prediction and forecast," in *Proceedings of the 18th ACM International Conference on Multimedia (MM '10)*, pp. 1235–1244, October 2010.

[23] M. Park, J. Luo, R. T. Collins, and Y. Liu, "Beyond GPS: determining the camera viewing direction of a geotagged image," in *Proceedings of the 18th ACM International Conference on Multimedia (MM '10)*, pp. 631–634, October 2010.

[24] L. Cao, J. Luo, A. Gallagher, X. Jin, J. Han, and T. S. Huang, "A worldwide tourism recommendation system based on geo-tagged web photos," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '10)*, pp. 2274–2277, March 2010.

[25] E. Gilbert and K. Karahalios, "Predicting tie strength with social media," in *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems*, pp. 211–220, April 2009.

[26] G. Yu, Y. Gu, and Y. B. Bao, "Large scale graph data processing on cloud computing environments," *Chinese Journal of Computers*, vol. 34, no. 10, pp. 1753–1767, 2011.

[27] R. Van Der Hofstad, *Random Evolution in Massive Graphs*, 2013.

[28] R. Dai and I. F. Akyildiz, "A spatial correlation model for visual information in wireless multimedia sensor networks," *IEEE Transactions on Multimedia*, vol. 11, no. 6, pp. 1148–1159, 2009.

[29] F. Jiangfan and S. Hu, "A data coding method of multimedia GIS in limited resource of mobile terminal," *Journal of Information & Computational Science*, vol. 9, no. 18, pp. 5873–5880, 2012.

[30] T. L. Le, M. Thonnat, A. Boucher, and F. Brémond, "Surveillance video indexing and retrieval using object features and semantic events," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 23, no. 7, pp. 1439–1476, 2009.

[31] T. Zhao and R. Nevatia, "Tracking multiple humans in crowded environment," in *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR 2004*, vol. 2, pp. II406–II413, July 2004.

[32] X. Y. Zhang, X. J. Wu, X. Zhou, X. G. Wang, and Y. Y. Zhang, "Automatic detection and tracking of maneuverable birds in videos," in *Proceedings of the International Conference on Computational Intelligence and Security (CIS '08)*, vol. 1, pp. 185–189, December 2008.

International Journal of

# Rotating
# Machinery

Journal of
Engineering

The Scientific
World Journal

Journal of
Sensors

International Journal of
Distributed
Sensor Networks

Advances in
Civil Engineering

Journal of
Control Science
and Engineering

Journal of
Robotics

Journal of
Electrical and Computer
Engineering

Submit your manuscripts at
http://www.hindawi.com

Advances in
OptoElectronics

VLSI Design

International Journal of
Navigation and
Observation

Modelling &
Simulation
in Engineering

International Journal of
Aerospace
Engineering

International Journal of
Chemical Engineering

International Journal of
Antennas and
Propagation

Active and Passive
Electronic Components

Shock and Vibration

Advances in
Acoustics and Vibration