# Context-Based Semantic Labeling of Human-Vehicle Interactions in Persistent Surveillance Systems

**Vinayak Elangovan**
vinayake@gmail.com
Dept. of Mech. & Mfg Engr.
Tennessee State University
TN, U.S.A.

**Amir Shirkhodaie**
ashirkhodaie@tnstate.edu
Dept. of Mech. & Mfg Engr.
Tennessee State University
TN, U.S.A.

**Abstract –** The improved Situational awareness in Persistent Surveillance Systems (PSS) is an ongoing research effort of the Department of Defense. Most PSS generate huge volume of raw data and they heavily rely on human operators to interpret and inference data in order to detect potential threats. Many outdoor apprehensive activities involve vehicles as their primary source of transportation to and from the scene where a plot is executed. Vehicles are employed to bring in and take out ammunitions, supplies, and personnel. Vehicles are also used as a disguise, hide-out, a meeting place to execute threat plots. Analysis of the Human-Vehicle Interactions (HVI) helps us to identify cohesive patterns of activities representing potential threats. Identification of such patterns can significantly improve situational awareness in PSS. In our approach, image processing technique is used as the primary source of sensing modality. We use HVI taxonomy as a means for recognizing different types of HVI activities. HVI taxonomy may comprise multiple threads of ontological patterns. By spatiotemporal linking of ontological patterns, a HVI pattern is hypothesized to pursue a potential threat situation. The proposed technique generates semantic messages describing ontology of HVI. This paper also discusses a vehicle zoning technique for HVI semantic labeling and demonstrates efficiency and effectiveness of the proposed technique for identifying HVI.

**Keywords**: Semantic Labeling, Human-vehicle interactions (HVI), Persistent Surveillance System (PSS), Zoning of Vehicle (ZoV), Soft Data (SD), Hard Data (HD)

## 1. INTRODUCTION

Multi-source information fusion is critical to the delivery of effective decision making in time critical applications. To realize the goals of interoperability and shared situation awareness in persistent surveillance systems, it has been recognized that fusion of information from soft (human) and hard (physical) sensors are the key enabling technology. The idea of fusing data from multiple sources is a compelling one, driven by the pursuit of a higher level of understanding that enables more effective action, which may be derived from the whole versus the sum of the parts. Yet, it is a challenging idea because the technologies for effective utilization of separate sources are not necessarily appropriate or adequate for discovering value and intelligence from the fused data. However, optimization of this process in the context of multi-modality homogeneous and heterogeneous sensors combining data from both hard and soft agent is a considerable challenge.

Hard Data (HD) i.e. information computed from hard sensors usually provides straight forward solution to analyze; their stability and characteristics are predictable [1]. HD requires low or moderate computing time to analyze the situation. Whereas Soft Data (SD) i.e. information computed from soft sensors (generally human resources), is strongly based on human intuition and the subjectivity. It represents the ambiguity in human thinking with real life uncertainty. Messages from human's i.e. SD sometimes provide valuable information or observations which may not be available from HD like in judging relationships between two individuals, their facial expressions etc. Information from soft sensors can be refined through various methodologies like Fuzzy Logic, Neural Networks, and Genetic Algorithms to generate well delineated semantic messages. Information extraction from both HD and SD is needed before any fusion can take place [2]. Semantic messages generated from hard and soft data are fused together to raise the percentage of certainty in analyzing the severity. Several issues are involved in fusing HD as in extracting the information from humans, quantifying the degree of certainty from human data, in modeling the human information and so on [3]. For example, the Times Square Car bomb detection incident on May 1$^{st}$, 2010 in New York had shown the importance of

SD. On the other hand, it also had increased fraught false alarms. For instance, a suspicious package alarmed by SD, turned out to be cooler filled with water bottles. This demonstrates how sorting out real time threat information from spoofs or from over reactive humans would be a burden of waste of time of investigators and cops who could have done more effective work [4].

To facilitate efficient fusion of data/information from physical sensors and soft (human) agents a common structure will suffice to ease integration of processed information into new knowledge. In this paper, we presented an approach for context-based semantic annotation of HVI. Here on, by HVI we refer the type of activities that a target human may exhibit while using his/her vehicle (example: opening/closing vehicle doors, hood or trunk, turning on/off engine, arriving/departing at/from a vehicle parking location). The objective of this paper is bi-folded. The first objective is to demonstrate that improved shared situational awareness can be achieved by fusion of data from multi-modality physical sensors (i.e., in our experiment we applied acoustic sensors to record sensors as they are generated while target human interacts in and around his/her vehicle). The second objective is to demonstrate automatic generation of informative semantic annotations facilitating fusion of discrete yet spatiotemporal correlated events under uncertainties. The major focus of this paper is, however, toward visual characterization of activities of a target human with a vehicle. A methodology is proposed for persistent tracking of target human activities around a vehicle including an approach for automatic generation of semantic annotations consistent with the behavior of the target human with the vehicle. The ultimate goal is to develop key techniques maximizing the information entropy and improving situational awareness a Persistent Surveillance System (PSS).

This paper is organized as follows: Introduction, Related Work, Methodology and Theoretical Background, Experiment Work and Result analysis, Conclusions and Future Work followed by Acknowledgements and References.

## 2. RELATED WORK

Several research activities have been done by previous researchers to improve the military intelligence using live video camera for fighting against crime and dangerous acts. Currently in many applications for detecting multiple suspicious activities through a real-time video feed, multiple analysts manually watch the live video stream to infer the necessary information [5]. Vision analysis is expensive and also prone to errors since humans have some limitations in monitoring continuous signals [6]. Video surveillance and Human behavior analysis is one of the ongoing active researches in image processing [5, 6, 7]. Other research activities focused on identifying spatiotemporal relations between human and vehicles for situation awareness.

Semantic concept detection has been explored for advancing machine learning. Supervised Learning methods like Neural Networks, Support Vector Machines, Naïve Bayes decision tree etc have been used effectively for detecting a number of visual concepts [8]. Detection of target of interest and image segmentation is significantly related tasks and can be used to detect events by feeding information from one task to other [9]. Image segmentation is partition of an image into different meaningful regions using similarities such as intensity, color, textures, patterns etc. But in real time situations, issues like limited spatial resolution, noise, poor contrast features, overlapping intensities etc makes segmentation a difficult task [10]. To overcome this issue, we had proposed a method of detection and tracking of target and events using target zoning.

Silhouette-based posture analysis was proposed in [11] to estimate the posture of the detected person. To identify whether a person is carrying an object in upright-standing posture, dynamic periodic motion analysis and symmetry analysis are applied in [11]. Events can be detected by imitating human - vehicle behavior through semantic hierarchy implementing Coordinate, Behavior and Event class, determining the behavior of individual vehicle in detecting the traffic events on the road [12].

In this paper, we have proposed a new approach of identifying and detecting HVI through Zoning of the target Vehicle. Initially, surrounding of a target vehicle is zoned into sub-regions for tracking human activities nearby the vehicle. Presence of human(s) within these spatial zones triggers events. To recognize a taxonomy of HVI, we inference an array of such spatiotemporal events to possible ontological patterns that leads to generating context-based semantic messages consecutively as events evolve over time. In our approach HVI taxonomy represents a hypothesis and our system can track multiple hypotheses on different threads.

## 3. METHODOLOGY AND THEORETICAL BACKGROUND

This section describes the zoning of target i.e. target being the vehicle; human and vehicles interactions; human, object and vehicles interactions; multi human and vehicles interactions, ontology development and message mapping followed by hard and soft fusion challenges.
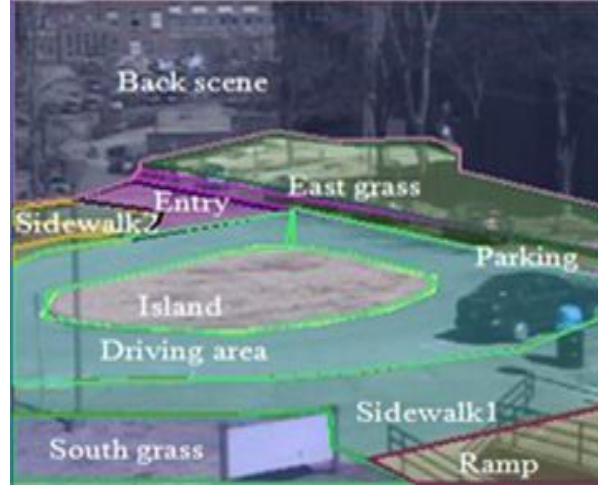


Figure-1: PSS Environment.

In order to generate context-based semantic labels in a PSS (as shown in Figure-1), we divide the monitored area into connected zones and label each zone with a label.

The conceptual architecture flow of the proposed system for detecting and semantically labeling the HVI is presented in Figure-2. The contextual information of the environment is used to semantically label geo-location of target vehicles. The heading direction of the target vehicle is determined as a vector as soon as it arrives into the surveillance environment. Zoning of Vehicle (ZoV), as described in next section, is employed for labeling pose of vehicle as it comes to a full stop. HVI activities around the vehicle are vigorously tracked via image processing techniques. Initially the environmental background is removed via a frame differencing technique. Next, a blob analysis and adaptive thresholding are performed to localize the human target around the vehicle. This results a binary image blob representing the human target. Next, the parameters of the blob are computed. These parameters include: length, width, area, central moments, and orientational angles. Blob of car have specific range of area. The final step is that of classifying structural posture of the target human as described in details in [18].

Each instance of such activity is considered as an event. Events are spatiotemporal. We generate a corresponding semantic message for each detected event. These generated semantic messages are matched against the known HVI Ontologies for adaptive scene semantic labeling. We have developed more than 200 ontological rules for describing HVI for any possible single vehicle - single human scenarios. These ontologies are maintained in a HVI Ontology library.
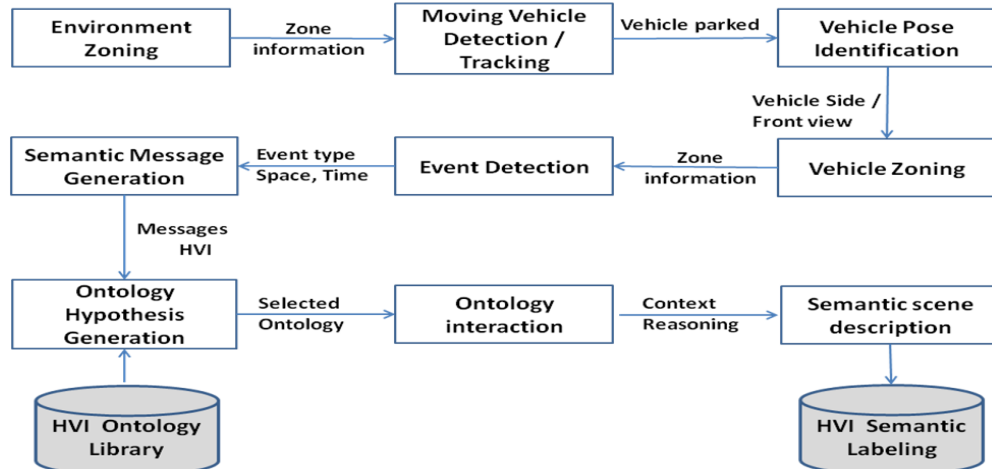


Figure-2: Framework of Semantic Labeling of HVI

### 3.1 Vehicle Spatial Zoning

By zoning surrounding of the vehicle and detecting human activities at such zones, one can ascertain type of potential/possible interactions that the human is involved with some degree of certainty. As the vehicle arrives and comes to a rest position, zoning is done for the vehicle. Two vehicle orientations (side view and frontal view) are illustrated in Figure 3. As shown, we partition the surrounding of vehicle into 20 different zones. The vehicle zoning helps in identying 'where', 'what' and 'when' an event had occurred around the vehicle. For example if a person opens the car hood, it can be identifed as event of "Hood Open" occurred in hood zone (i.e. zone-15) and semantic messages are generated accordingly to describe the HVI. As mentioned earlier, semantic labeling of events is generated with certain degree of confidence. To reduce false alarm rate, we fuse information from two or more views of the vehicle (when available) and apply a bayesian belief approach to further improve signal to noise content and reduce uncertainty associated with characterization of HVI events.
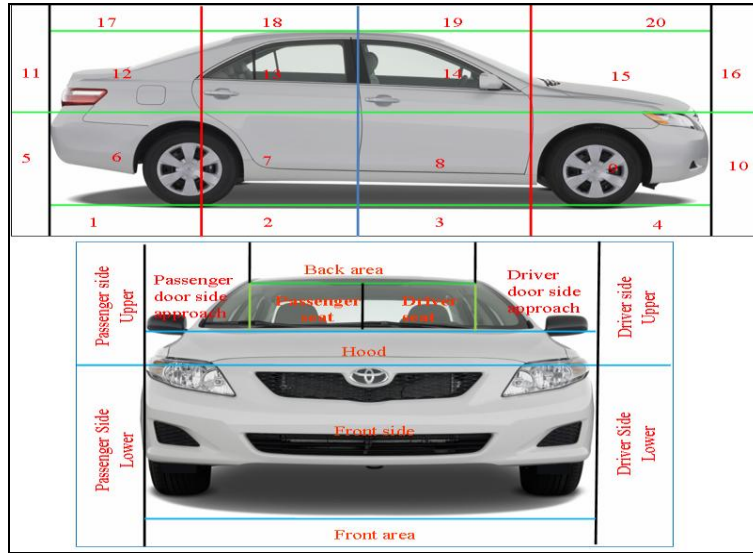


Figure 3. Top - Zoning of Vehicle Target in Side View, Bottom - Zoning of Vehile Target in Front View

### 3.2 Taxonomy of Human-Vehicle Interactions

The taxonomy of HVI is presented in Figure 4. The taxonomy has a hierarchical structure and presents three types of interactions that a human and a vehicle may exhibit. Categorically, such interactions include: (1) visual identification, (2) human action identification, and (3) auditory identification. Each category is further branched out to define more atomic HVI relationships, not shown here in brevity of space limitation.

### 3.3 Characterization and Semantic Labeling of HVI

A sequence of activities involving human-vehicle interaction with a parked vehicle is shown in Figure 5. The activities include: Image processing is performed for identifying the blobs of interest and categorizing events. For example in (f), it is identified an event had occurred in hood zone i.e. hood open and human is not found in any of the car zones which triggers a message as *'A human had opened car hood and left the scene'*.
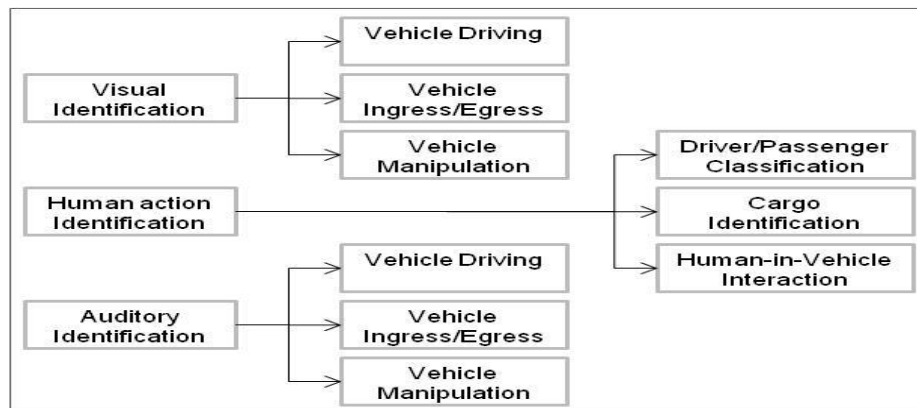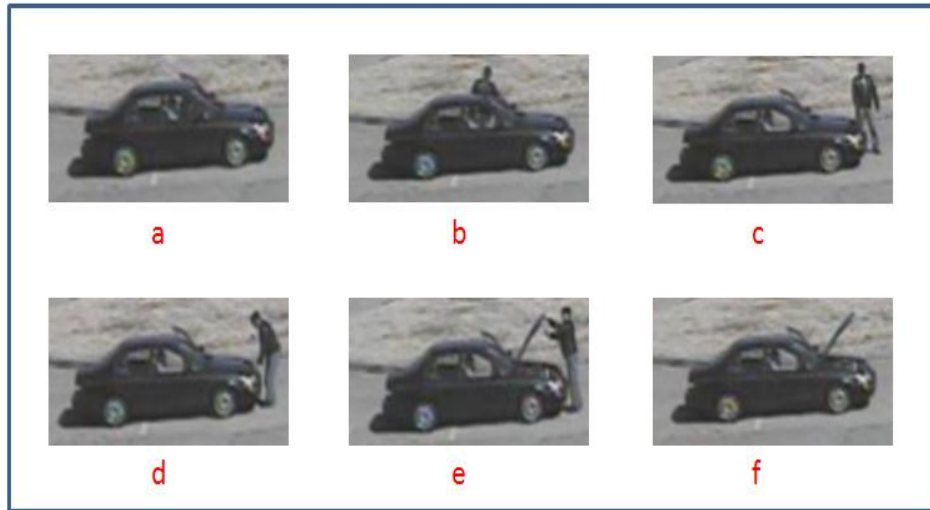


Figure 4. Taxonomy of HVI

Figure 5. (a) – driver opens driver-side door, (b) – driver gets off from the driver side door, (c) – driver walks towards the car hood, (d) – drive attempts to open up the car hood, (e) – driver opens the car hood, (f) – driver leaves the vehicle after opening the car hood.
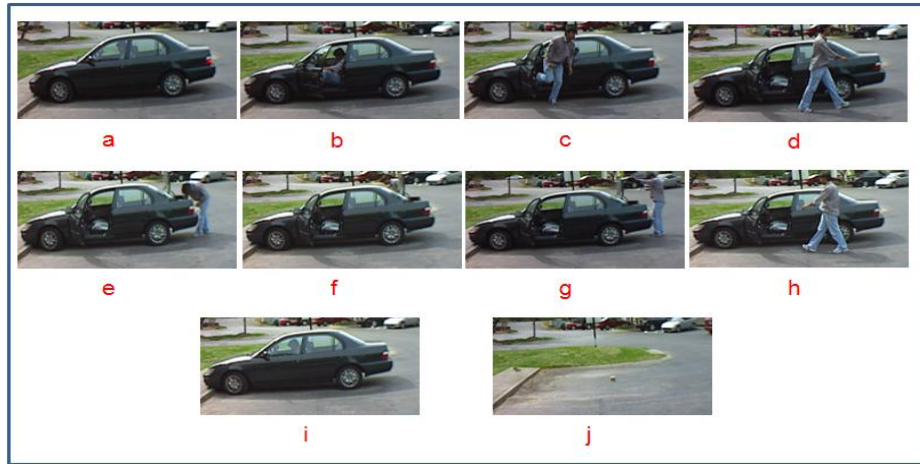


Figure 6. Senario-2 - Sequence of Human-Object-Vehicle Interaction. (a) a vehicle stops at one the parking area, (b) driver opens the driver side door, (c) – driver walks out from the driver side leaving the door open, (d) – driver walks towards the trunk, (e) driver opens the trunk and takes a small object, (f) – driver drops the object behind the back passenger side door, (g) – driver closes the hood, (h) driver walks towards the driver side door, (i) – driver get into the driver's seat; and (j) – vehicle leaves while an unattended package is left behind in the scene.

Another HVI is illustrated in Figure 6. In this simulated example. The human target leaves the scene with a box without closing the car trunk still open. Figure-6 shows the sequence of events involved in human, object and vehicle interaction in Scenario-2. As discussed in section 3.3, our technique is able to track multi human-vehicle interactions using the proposed methodology. Figure-7 illustrates one such example where (a) target vehicle stops in a restricted zone, (b) driver side doors and south side back seat door open up, (c) – two persons coming out from driver side doors, (d) – third person comes out, (e) vehicle three doors are found open, finally (f) – an individual is found near the driver side passenger door.
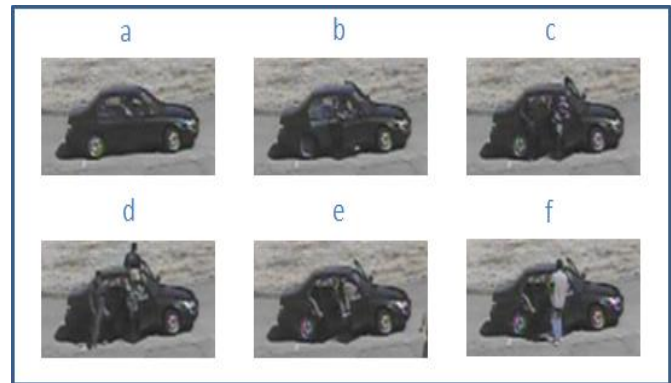


Figure 7. Scenario-3 - Sequence of Multi-Human-Vehicle Interaction.

We use Resource Description Framework (RDF) Prime for generating semantic scene labels. The RDF is a language for representing information about resources in the World Wide Web. RDF has features that facilitate data merging even if the underlying schemas differ, and it specifically supports the evolution of schemas over time without requiring all the data to be changed. The semantic messages we generate has three main parts: "Subject", "Predicate", and "Object" and two augmented information that present the time and place where the event has occurred. 'Predicate' contains an action verb and may contain objects, adjectives and adverbs. Each generated semantic label is followed by the time of event had occurred. A complete generated HVI semantic messages, therefore has the following protocol format: '|Subject|, |Predicate|, |Object|, |Time|, and |Place|'. For example, *'Driver Opened the Vehicle Hood@11:40:18#ParkingLot'* or *'Driver Lifted up the vehicle Trunk @08:23:24#Zone21'*.

### 3.4 Human-Vehicle Interaction Ontologies

Many disciplines now develop standardized ontologies that domain experts can use to share and annotate information in their fields. By definition, the ontology is explicit formal specifications of the terms in the domain and relations among them [19]. In another words, the ontology is involved with an iterative method of Knowledge-Engineering (KE) for a specific domain.

The HVI Ontology development has several advantages including: (1) it facilitates common sharing of situational awareness, (2) it enables reuse of domain knowledge, (3) it makes domain assumptions explicit, (4) it separates domain knowledge from the operational knowledge, and (3) It facilitate to analyze domain knowledge. By focusing on metaphysics of HVI, we develop a rule based system containing 200+ rules that links together what types of HVI are possible and what relations these events bear to one another to ensemble a situational awareness. Figure 8 illustrates our hierarchical structure of the HVI ontologies.
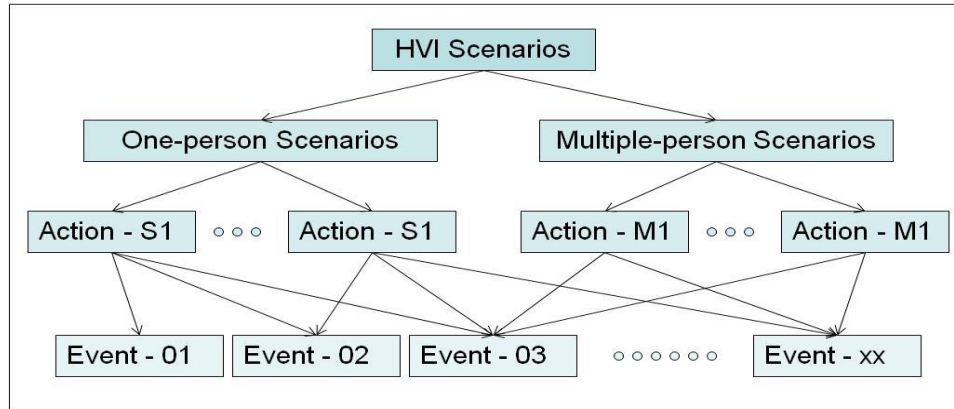


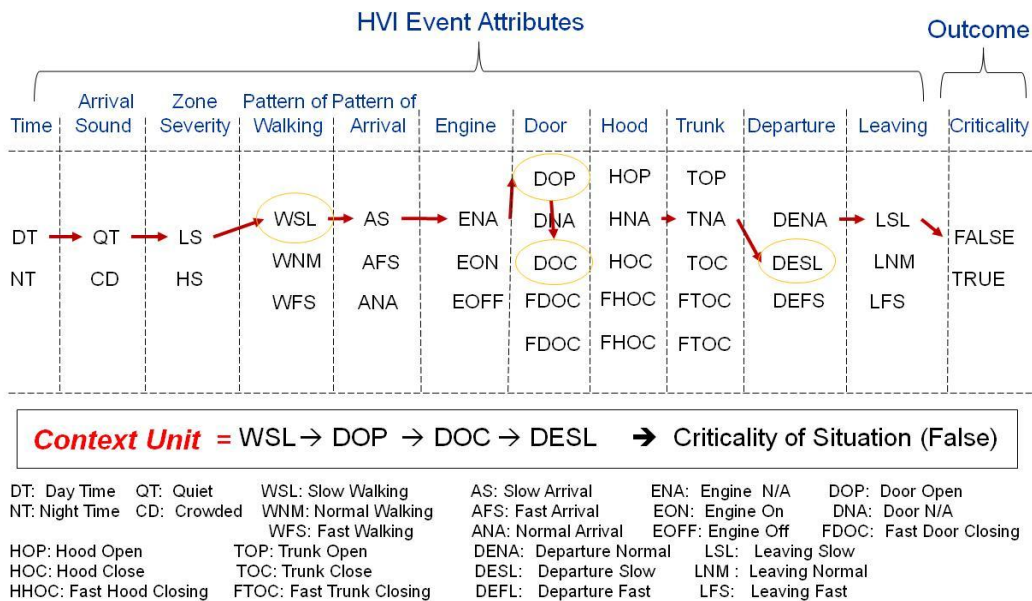Figure 8. Illustrates the Hierarchy of HVI Ontologies.



Figure 9. HVI events properties.

Figure 9 presents the breakdown of atomic events and their compositional "Context Unit" that represents a specific HVI scenario along with an assigned HVI situational criticality level. As mentioned earlier, the HVI ontology development is a form of knowledge engineering. By means of this ontological approach, we selectively characterize unique HVI context unit. As demonstrated in Figures 10 and 11, our HVI ontology is presented in tree structure. The ontology tree is constructed based on clustering of ordered atomic events. List of atomic events in development of the HVI ontology are listed at the bottom of Figure 8. One main advantage for presenting the HVI ontologies in the tree structure is that more complex ontologies can be developed based on simpler ontologies much more efficiently.
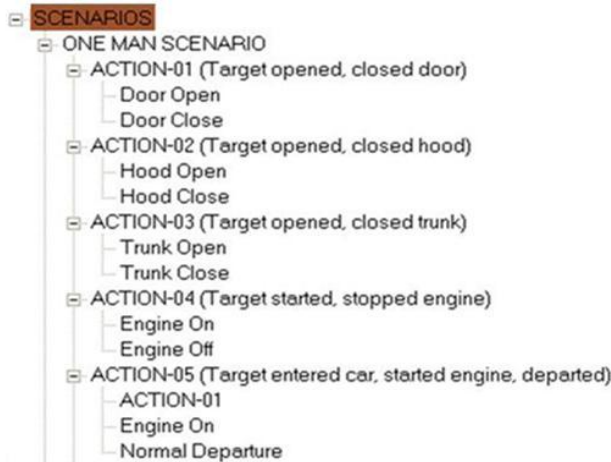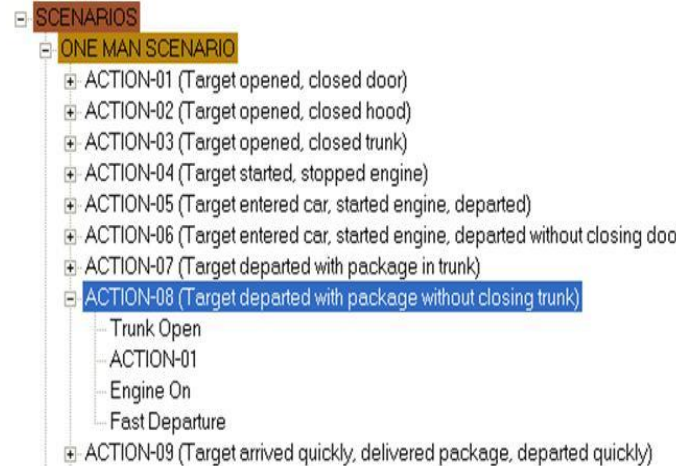


Figure-10: HVI Ontology in Tree Structure



Figure 11. More Complex HVI Ontologies.

The interaction of humans with car can have a large number of scenarios, particularly, depending on time, space, and environmental factors criticality level of each ontology thread may vary. While discussion of such event modulation factors is outside the scope of this paper. The interested reader, are encouraged to refer to our other publication paper [13] for further details on this discussion. The framework of our soft and hard sensor data fusion process is presented in Figure 12. This framework accommodates integration of information from both soft agents (i.e., human observers reporting events
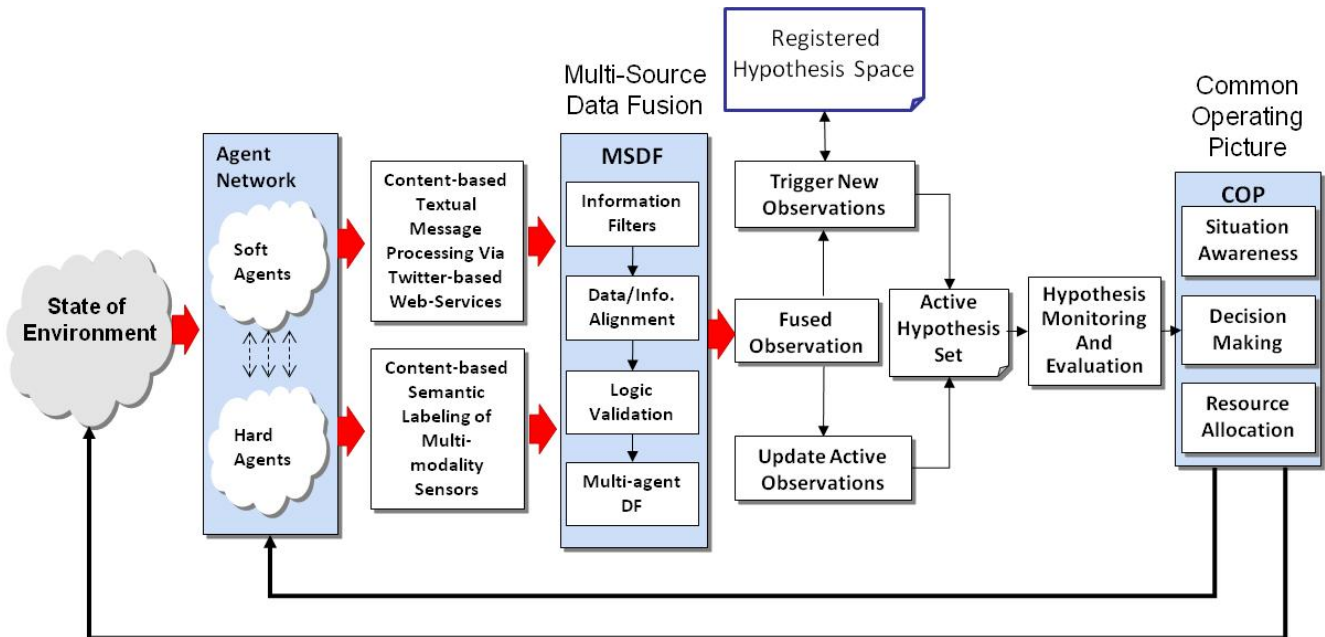


Figure 12. A Framework for Soft-Hard Sensor Data Fusion.

(e.g., via a twitter-bsed web-service [14])) as well as hard agents (i.e., physical sensors whose information are modulated via a context-based environment inferencing [13]).

## 4.0 EXPERIMENTAL WORK

This section describes an experiment carried out for validation of semantic annotation of human-vehicle interactions in a context of a simulated scenario as depicted in the Figure 5. In this simple scenario, the driver parks the vehicle, opens up the trunk and removes an object from the trunk and leave the object on opposite site of the vehicle and returns back to driving seat and leaves the scene. Though nothing peculiarly critical about this scenario, we intentionally set up this experiment to verify how well our ontology-based approach can assist in generating meaningful and appropriate semantic messaging annotating different stages of interactions of the driver with his vehicle. Figure 12 illustrates an overview of our experiment. In this experiment we used one camera and one acoustic sensor. The semantic labeling of acoustic signatures is not discussed here and a complete discussion on that can be found in our other paper at this conference [20].

The seven main steps of this image processing process included:

1. Image frame differencing for isolating vehicle
2. Binary vehicle and obtain its edge silhoutte
3. Use vehicle silhoutte as a template for segmenting human from background (i.e., vehicle)
4. Use Zoning of vehicle technique to localize human around the vehicle
5. Human gating and posture Tracking
6. Human posture classification
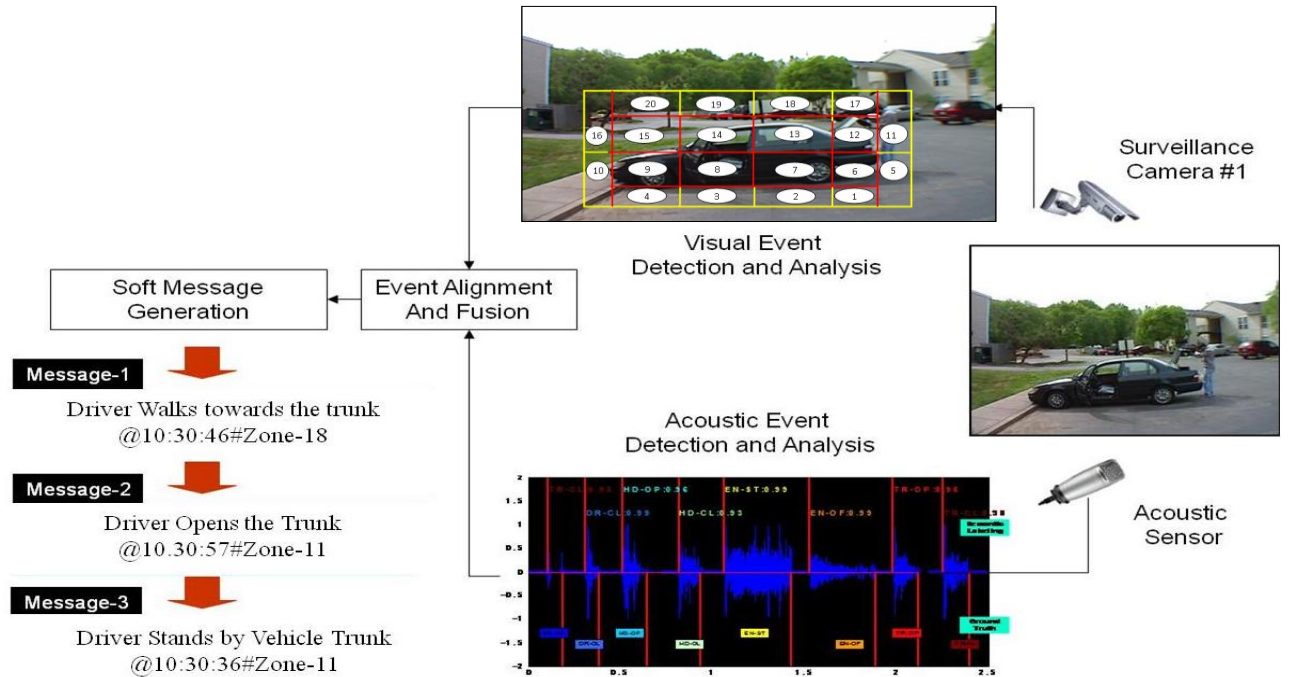7. Generate semantic annotation of HVI



Figure 13. Illustrates an overview of our experimental setup and different sensors used for semantic annotation of HVIs.

Figure 13 illustrates some of image processing steps for segmentation of vehicle and human through this scenario. Our image procesing technique generates ontology-based HVI semantic annotations per event as described in section 3.4. Each semantic annotation follows the specific protocol that is:

**HVI Protocol:** [Subject-Predicate-Object-Time-Space]

The *Subject* identifies the human target of interest. The *Predicate* defines an action verb describing the state of human target. The *Object* may take different forms. It can be either some material that may be perceived by the senses, or a noun or noun equivalent, the goal or end of an effort or activity. The *Time* tag aids in ordering and synchronizing randomly arriving semantic annotations. The *Space* tag ensures the semantic annotations are related to the same state of affairs in the environment. Figure 14 shows the steps of Image Processing in Human and Vehicle Segmentation in tracing HVI activities. Figure 15 presents the semantic annotations of the conducted experiment.
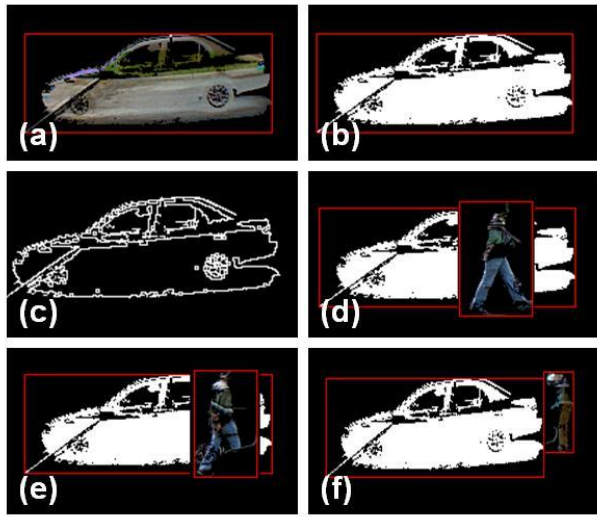


Figure 14. Image Processing Steps for Human and Vehicle Segmentation. (a) Vehicle Profile via Frame Differencing, (b) Binary Vehicle, (c) Vehicle Silhoutte, (d) Driver segmentation from background vehicle, (e) same as (d), and (f) human opening the trunk.

## 5. CONCLUSIONS AND FUTURE WORK

This paper presented a technique for semantic annotation of context based human-vehicle interactions. A vehicle zoning technique is also proposed for interaction detection and tracking of human(s) around the vehicle. Developed system has been tested for one human and one vehicle interaction as well as multiple human and one vehicle. In our experiment two sensor modalities (i.e., imaging and acoustic) are considered. For each sensor modality, we generate consistent and complementary semantic annotations. The inferred sensor information can be fused readily to describe evolving HVI's in an evocative way. Lastly we have proposed a new framework for integration of soft data and hard data and for generating further carry out fusion of soft and hard data to raise the degree of confidence and in describing the risk of certainty involved in suspicious HVI activities.
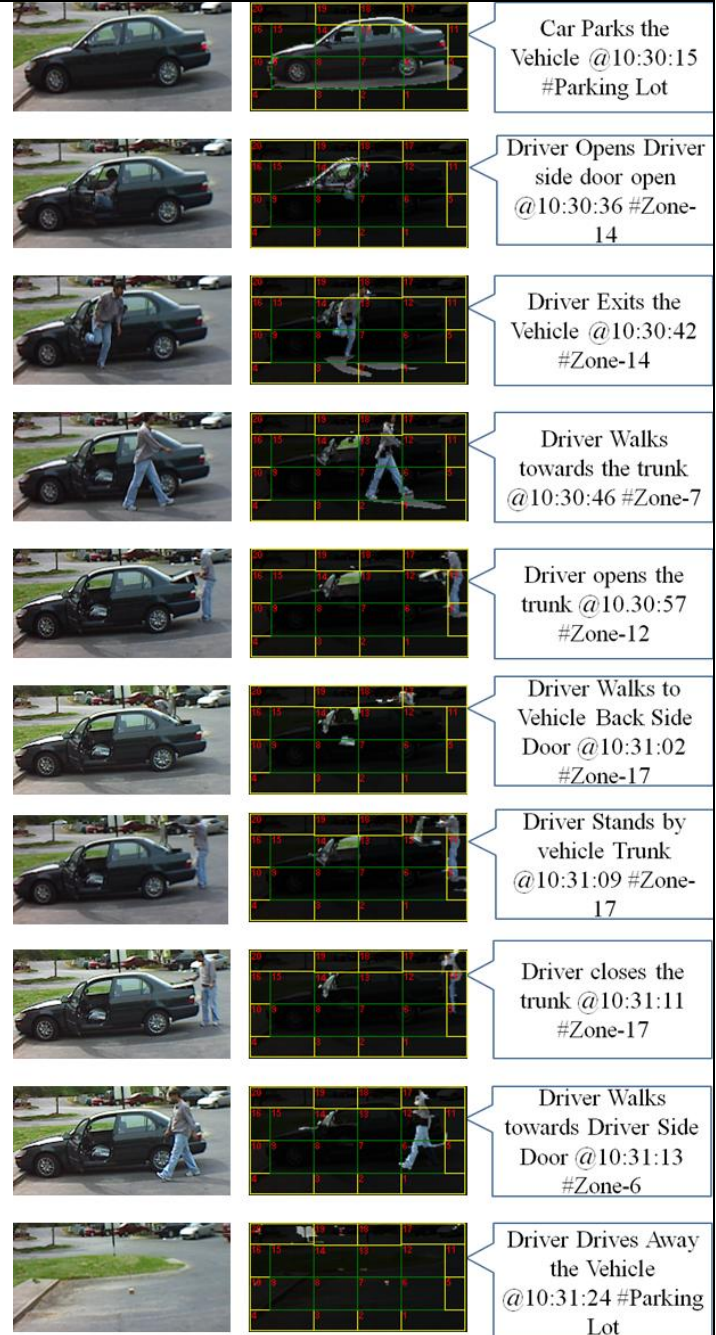


Car Parks the Vehicle @10:30:15 #Parking Lot

Driver Opens Driver side door open @10:30:36 #Zone-14

Driver Exits the Vehicle @10:30:42 #Zone-14

Driver Walks towards the trunk @10:30:46 #Zone-7

Driver opens the trunk @10.30:57 #Zone-12

Driver Walks to Vehicle Back Side Door @10:31:02 #Zone-17

Driver Stands by vehicle Trunk @10:31:09 #Zone-17

Driver closes the trunk @10:31:11 #Zone-17

Driver Walks towards Driver Side Door @10:31:13 #Zone-6

Driver Drives Away the Vehicle @10:31:24 #Parking Lot

Figure 15. Detection of events in HVI and Semantic message generation

## ACKNOWLEDGMENTS

## REFERENCES

[1] Ovaska, S.J., VanLandingham, H.F. and Kamiya, A., "Fusion of soft computing and hard computing in industrial applications: an overview," IEEE Trans. Syst. Man Cybernet. Part C: Appl. Rev. v32. 72-79.

[2] Marco A. Pravia, Olga Babko-Malaya, Michael K. Schneider, James V. White, and Chee-Yee Chong, "Lessons Learned in the Creation of a Data Set for Hard/Soft Information Fusion," 12th International Conference on Information Fusion, July 6th – 9th, 2009, Seattle, WA.

[3] D.L. Hall, J. Llinas, M. McNeese, and T. Mullen, "A Framework for Dynamic Hard/Soft Fusion," Proc. 11th Int. Conf. on Information Fusion, Cologne, Germany, July 2008.

[4] http://blog.newsweek.com/blogs/declassified/archive/ 2010/05/07/what-can-intelligence-agencies-do-to-spot-threats-like-the-times-square-bomber.aspx.

[5] Joshua Candamo, Matthew Shreve, Dmitry B, Goldgof, Deborah B. Sapper, and Rangachar Kasturi, "Understand Transit Scenes: A Survey on Human Behavior – Recognition Algorithms" in IEEE Transactions on Intelligent Transportation Systems, Vol.11, No.1, March 2010.

[6] N. Sulman, T. Sanocki, D. Goldgof, and R. Kasturi, "How effective is human video surveillance performance?" in Proc. Int. Conf. Pattern Recog., 2008, pp. 1–3.

[7] W. Hu, T. Tan, L. Wang, and S. Maybank, "A survey on visual surveillance of object motion and behaviors," IEEE Trans. Syst., Man, Cybern.C, Appl. Rev., vol. 34, no. 3, pp. 334–352, Aug. 2004.

[8] Shih-Fu Chang, Wei-Ying Ma, Arnold Smeulders, "Recent Advances and Challenges of Semantic Image/Video Search," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Hawaii, USA, April 2007.

[9] Stephen Gould, Tianshi Gao, Daphne Koller, "Region-based Segmentation and Object Detection," Proceedings of Advances in Neural Information Processing Systems (NIPS), 2009.

[10] Yong Yang, "Image segmentation based on fussy clustering with neighborhood information," Optica Applicata, Vol. XXXIX, No.1, 2009.

[11] Ismail Haritaoglu, David Harwood and Larry S. Davis, "W4: Real-Time Surveillance of People and Their Activities," IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.22, No.8, August 2000.

[12] Shunsuke Kamijo, Masahiro Harada and Masao Sakauchi, "An Incident Detection System Based on Semantic Hierarchy," IEEE Intelligent Transportation Systems Conf., Oct. 3rd–6th, 2004, Washington, DC.

[13] Shirkhodaie, A., "Rababaah, H. " Multi-Layered Impact Modulation for Context-based Persistent Surveillance Systems, " in the SPIE 2010 Defense, Security, and Sensing Conference, Multi-Sensor, Multi-Source Information Fusion: Architectures, Algorithms, and Application, April 2010, Orlando, FL.

[14] Rababaah H., and Shirkhodaie, A., "Twitter-based Web-service for Human Observation Aggregation in Hybrid Sensor Networks," in the SPIE 2009 Defense, Security, and Sensing Conference, Multi-Sensor, Multi-Source Information Fusion: Architectures, Algorithms, and Application, April 2010, Orlando, FL.

[15] Rababaah, H., and Shirkhodaie, A., "Soft Adaptive Fusion of Sensor Energy (SAFE) in Large-Scale Sensor Networks, " SPIE Defense and Security Conf., Multi-Sensor, Multi-Source Information Fusion: Architectures, Algorithms, and Application, paper 7345-07, April 13-17, 2009, Orlando, FL.

[16] Rababaah, H., and Shirkhodaie, A., "Feature Energy Assessment Map (FEAM): a Novel Model for Multi-Modality Multi-Agent Information Fusion in Large-Scale Intelligent Surveillance Systems, Networks, " SPIE 2009 Defense, Security, and Sensing Conference, Multi-Sensor, Multi-Source Information Fusion: Architectures, Algorithms, and Application, paper 7345-19, April 13-17, 2009, Orlando, FL.

[17] Rababaah H., and Shirkhodaie A., "A Survey of Intelligent Visual Sensors for Surveillance Applications," the IEEE Sensor Application Symposium, February 6-8, 2007.

[18] Rababaah H., Shirkhodaie A., "Human Posture Classification for Intelligent Visual Surveillance Systems, " in the SPIE Defense and Security, Visual Analytics for Homeland Defense and Security, March 17-20, 2008, Orlando, FL.

[19] Xiquan Yang, Ye Zhang, Na Sun, Deran Kong, "Research on Method of Concept Similarity Based on Ontology," Proc. of the 2009 International Symposium on Web Information Systems and Applications (WISA'09), *China, May 22-24, 2009, pp. 132-135.*

[20] Rababaah, H., Shirkhodaie, A., "Semantic Labeling of Non-stationary Acoustic Events in Persistent Surveillance System," SPIE Defense, Security and Sensing, April 2011, Orlando, Florida.