

Reliable Extraction of Semantic Information and Rate of Innovation Estimation for Graph Signals

Mert Kalfa, Sadık Yağız Yetim, Arda Atalik, Mehmetcan Gök, Yiqun Ge, Rong Li, Wen Tong, Tolga M. Duman, and Orhan Arıkan

Abstract—Semantic signal processing and communications are poised to play a central part in developing the next generation of sensor devices and networks. A crucial component of a semantic system is the extraction of semantic signals from the raw input signals, which has become increasingly tractable with the recent advances in machine learning (ML) and artificial intelligence (AI) techniques. The accurate extraction of semantic signals using the aforementioned ML and AI methods, and the detection of semantic innovation for scheduling transmission and/or storage events are critical tasks for reliable semantic signal processing and communications. In this work, we propose a reliable semantic information extraction framework based on our previous work on semantic signal representations in a hierarchical graph-based structure. The proposed framework includes a time integration method to increase fidelity of ML outputs in a class-aware manner, a graph-edit-distance based metric to detect innovation events at the graph-level and filter out sporadic errors, and a Hidden Markov Model (HMM) to produce smooth and reliable graph signals. The proposed methods within the framework are demonstrated individually and collectively through simulations and case studies based on real-world computer vision examples.

Index Terms—Semantic signal processing, semantic communications, semantic graph signals, goal-oriented signal processing, goal-oriented communications.

I. INTRODUCTION

RECENT advances in machine learning techniques are resulting in a paradigm shift in signal processing, where semantically-rich information about any underlying signal is becoming increasingly available. Extraction of semantic information from videos [1], sounds [2], [3], financial time-series data [4], and many more signal modalities enables processing, storage, and communication at a semantic level. As a result, the next generation of signal processing and communication systems are expected to include *semantically-aware* agents that can optimize the processing and transmission of information over the inherent semantic meaning [5]–[8].

In [9], we introduced a goal-oriented signal processing framework, where any input signal can be mapped to an application-specific graph-based semantic language, in conjunction with internal or external goals that are also based on the same language. Compared to the natural-language-based applications of the semantic processing techniques, the application-specific nature of the proposed semantic language in [9] reduces the computational requirements of the agents at the site of the sensors, and provides inherent privacy and secrecy by mapping raw data into graph-signals.

In the goal-oriented semantic signal processing framework [9] a *semantic extractor*, which is typically a machine learning algorithm, transforms the input signal into graph

signals with a hierarchical structure that also include embedded numerical attributes. Regardless of the adopted semantic signal processing framework, this transformation process in any semantics-enabled system will produce noisy semantic signals due to the signal-to-noise-ratio (SNR) of the input signal (at the technical level), insufficient training, model imperfections, and knowledge-base limitations. For an object classifier implemented on a video signal, the low input SNR can be due to the distance of the objects or due to the improper lighting conditions, whereas model errors can be due to lack of training or simply the limitations of the model itself (i.e., number of layers, structure of the neural network, number of neurons, etc.). The knowledge-base for the framework in [9] can be defined as the range of outcomes of the pre-defined language structure and the logical likelihoods of certain outcomes (see Section III for details), where a limited definition can lead to missed detections or wrong classifications. Note that the semantic noise discussed here is not the semantic channel noise [10], nor is it an adversarial agent distorting the input of the algorithm [11] but rather a source noise introduced when the semantic signal is first generated. An advantage in handling the source noise is that we can have access to the statistical characteristics of the semantic extractor, as well as logical characteristics of the semantic information, and exploit these to improve the fidelity of the results.

Goal-oriented filtering of the semantic information helps focus the remaining signal processing on those graph signals that are of interest. For both storage and transmission purposes, it is important to detect innovation events on the extracted semantic signals [9]. The proposed semantic signals in [9] have a hierarchical structure with graphs and numerical attributes; therefore, innovation events can happen at different levels of the semantic signal. An example of an innovation event at the graph signal level can mean a graph pattern of interest (that is predefined by the internally or externally defined goals) has emerged (or receded) in the semantic description of the raw signal. Once a graph pattern of interest is identified and is being tracked, the numerical attributes such as position and subfeature vectors can be tracked across time to detect innovation events at the attribute level (again, based on the interests defined by the goals). To give a more practical example, consider a computer vision application that is used to provide security for a pedestrian-only street. In this scenario, the external goals can be the detection of objects in the classes of *vehicles* and *suitcases*. Therefore, innovations at the graph level can be identified as events where these classes of interest come into our out of the semantic graph description of the scene. Additionally, once these classes are identified, their

accurate positions, as well as some of their features (a large suitcase is more attention-worthy than a small one, etc.) can be tracked across time to be stored or transmitted once they exceed a certain threshold.

Quantification of the innovation events enables efficient and accurate tracking of the semantic signals, filtering out sporadic erroneous detections, and scheduling the storage or transmission events, depending on the application. Coupled with the inherent compression provided by a semantic language [9], [12], a scheduler based on semantic innovations can reduce the overall transmission traffic dramatically.

In this paper, with the objective of improving the fidelity and quantifying the innovation of semantic information, we propose and demonstrate a semantic extraction framework that can be implemented on graph signals with embedded numerical attributes. Specifically, we introduce an integration technique to improve semantic fidelity in a class-aware manner, a modified graph edit distance (GED) metric and a Hidden Markov Model (HMM) to quantify, track, and smooth the semantic signals at the graph level, and a subspace tracking algorithm to quantify and track the semantic information at the numerical attribute level. The proposed framework can be implemented in the next generation of communication networks and sensor devices to enable semantic signal processing and semantic communications. The individual methods can be used either separately or collectively for filtering erroneous detections, fusion of multiple sensors/data-streams, and transmission scheduling for semantic communication applications. The methods presented in this paper are rigorously explained and demonstrated using simulations and computer vision examples on acquired video signals.

The rest of this paper is organized as follows. Section II provides a short literature review on the semantic signal processing, and semantic graph signal extraction. Section III briefly re-introduces the semantic signal processing framework of [9]. Section IV presents the proposed semantic extraction framework and its main building blocks. The details of each block in the proposed framework are given in Sections V–VIII. Section IX demonstrates the performance of the proposed methods via simulations and real-world computer vision case studies. A short discussion on the rate of innovation of semantic signals and the corresponding data throughput is given in Section X. Concluding remarks and future research directions are given in Section XI.

II. RELATED WORKS

The initial work on semantic information is almost as old as Shannon’s seminal work [13], with Weaver introducing the semantic communication paradigm [14], then Bar-Hillel and Carnap proposing a Semantic Information Theory (SIT) built on a semantic language based on propositional logic. In the last decade, the quantum leap in the abundance of semantically rich data enabled a resurgence of research on semantic signal processing and communications. To give a clear overview of the state-of-the-art on this emerging but highly active field, we review the literature in the following two subsections, namely on extraction of semantic information and on semantic communications.

A. Extraction of Semantic Information

Semantic information exists in many signal modalities such as a textual description of an image, a knowledge graph derived from a paragraph, and even in correlation functions of random processes. We refer to semantic transformation or semantic extraction as the mapping from an input modality to a target semantic modality where the target semantic modality can be anything ranging from vectors, texts, and graphs. Most popular semantic transformation techniques (especially in computer vision applications) include object detection [15]–[20], semantic segmentation [21]–[33], and captioning [34]–[47].

A powerful form of semantic transformation is to convert signals into graphs that represent a scene and encode the relationships presented in the signal. Scene graphs are proposed in [48] describe image features and object relationships in an explicit and structured way for image retrieval. Some recent papers [49]–[51] also consider joint optimization of object detection and relationship recognition parts. Specifically, Factorizable-Net is proposed in [52] where a Region Proposal Network (RPN) is used to extract object proposals and proposed objects are paired to obtain a fully-connected initial coarse graph. In [53], Graph-RCNN is introduced. Graph-RCNN uses relation-proposal-network (RePN) to prune the connections in the initial graph and an Attentional Graph Convolutional Network [54] refines the features on the graph. On the other hand, VCTree model [55] constructs a dynamic tree from a scoring matrix where visual context is encoded into the tree structure. Furthermore, some recent works [49], [50] use recurrent architectures for graph inference. Particularly, in [50] a feature refining module consisting of edge and node Gated Recurrent Units (GRU), and in [56], stacked bi-directional Long Short-Term Memory models (bi-LSTMs) are used.

Scene graphs can also be extracted from video signals. In [57], each frame in the video is converted into a scene graph as an intermediate semantic representation. Then, using frame and cross-frame level relationships of intermediate scene graphs, a story of the video is generated. Joint parsing of cross-view videos is introduced in [58] where scene-centric and view-centric graphs are hierarchically generated. For more details, we refer the readers to survey papers [59], [60].

B. Semantic Communications

Many outlook papers on next-generation communication networks envisioned that some form of semantic communications will pave the way for next generation wireless communications [6], [10], [61]–[66]. Recently more practical approaches for viable semantic communications are presented in the literature. In [12], a deep learning enabled semantic communication architecture called DeepSC is proposed for text transmission. Unlike traditional approaches where the objective is to minimize bit errors, this work focused on minimizing semantic errors while maximizing the system capacity with a mutual information estimation model for semantic channel coding. In [67], the authors extend their previous work to affordable internet-of-things (IoT) applications where

a lite and distributed version of [12] called the L-DeepSC is proposed. In [68], the authors propose a semantic communications framework where natural language texts are modeled as knowledge graphs that are transmitted and converted back into text at the destination. The conversion into and from knowledge graphs offer an intuitive and explainable semantic extraction compared to [12], [67].

Transmission of speech data for semantic communications is considered in [69]. The authors proposed a deep learning enabled system called DeepSC-S which extracts essential semantic information from speech audio signals by using squeeze-and-excitation (SE) networks [70] and applies attention based weighting to emphasize the essential information where source and channel encoder/decoder is jointly designed to mitigate channel distortion. The proposed architecture is shown to outperform traditional approaches in almost any SNR regime.

In [71], a semantic communication architecture (SC-AIT) is introduced where both effectiveness, semantic and technical level of communication is inter-connected via a neural network supported module. The proposed architecture is realized on a real-world test-bed for image classification task to detect surface defections. The authors showed that the proposed architecture outperforms traditional approaches in any SNR regime while having lower latency and higher compression ratio.

Real-time semantic communications is first studied in [72] where the authors developed a prototype for wireless image transmission based on the field-programmable gate array (FPGA), Vision Transformer (ViT) [73] and a denoising auto-encoder. The implemented prototype has been shown to be superior to traditional 256-quadrant amplitude modulation (256-QAM) in the low-SNR regime on a measure of structural similarity index (SSIM). In [74] constellation constrained version of DeepJSCC [75] (DeepJSCC-Q) is proposed for wireless image reconstruction. DeepJSCC-Q utilizes a differentiable soft quantization layer to map latent semantic vectors to transmitted symbols such that each quantization level corresponds to a learnable constellation point. It is illustrated that DeepJSCC-Q achieves comparable performances with respect to its unconstrained counterpart and outperforms the traditional methods which rely on separate source and channel coding scheme.

Deep learning aided end-to-end JSCC for wireless video transmission scheme (DeepWive) is introduced in [76] where the proposed architecture performs video compression, channel coding and modulation with a single neural network, resulting in direct mapping from video signals to channel symbols whereas the trained semantic decoder predicts the residuals without distortion feedback.

In [77], effectiveness of semantic communications is studied within a *joint learning and communication framework*. The authors adopted multi-agent reinforcement learning (MARL) approach to develop a systematic structure for collaborative agents participating in treasure hunts. The problem is formulated as a multi-agent partially observed Markov decision process (MA-POMDP). Agents are intended to learn policies to effectively exchange messages over a shared noisy wireless

channel for improved coordination and collaboration while taking long-term rewarding actions. As a result, the proposed joint learning and communication framework achieves higher performance than treating action-taking and communication aspects of MARL separately.

In the following section, a brief introduction to our work on semantic information and signal processing [9] will be presented.

III. PRIMER ON THE GOAL-ORIENTED SEMANTIC SIGNAL PROCESSING FRAMEWORK

Recently a goal-oriented semantic signal processing framework has been proposed to extract and process semantic information generated at sensor nodes [9]. In this framework, the semantic information is organized using a flexible graph-based language in a structured and hierarchical way. The highly organized and hierarchical approach to construct the semantic language shifts the processing load to the initial semantic information extraction, which in turn leads to ease of processing, filtering, and reduce the overall storage within a device or transfer of information between devices within a network.

The proposed structure is made up of bipartite graphs that include the identified signal *components*, and *predicates* that show the state or relationship of the detected components. Each node in the graphs may also include layers of numerical attributes with different levels of complexity. Although the aforementioned structure of the proposed language is fixed, the particular definitions of the components, predicates, and attributes are application-specific. Therefore, the structures are only as complex as the specific requirements of the implementation at the sensor/agent level and the specific situation at the sensor site. In comparison to the natural-language-based semantic signal processing applications [78], the proposed structure for the semantic information offers a much more efficient and unambiguous representation while inherently providing privacy and secrecy through relatively abstract graph signals. Moreover, the additional computational cost of introducing graph signal processing algorithms is balanced by the inherent low-order graph signals generated by the proposed framework as they are generated within an application-specific and limited dictionary.

The proposed structure of the goal-oriented semantic signal processing framework at a sensor node is shown in Fig. 1. Briefly, a raw input signal generated by the sensor is preprocessed (e.g., through up/downsampling or Fourier transform) before *semantic extraction*, where the sensor output is mapped to the target semantic output. Semantic filtering and post-processing blocks are implemented in a goal-oriented manner to reduce the range of semantic outputs to the specific goals and perform additional processing (e.g., source coding), respectively. Finally, the resulting goal-filtered and compressed semantic output is either transmitted or stored. The most critical component in the entire diagram shown in Fig. 1 is the semantic extraction block, since accurate and reliable mapping of the raw signals into the semantic language affects the performance of the remaining blocks that follow in the

semantic signal processing chain. Before we present the methods and algorithms to improve the accuracy and reliability of the semantic extractor, we first review the semantic language structure of [9].

As detailed in [9], the semantic description of a signal starts with the definition of component and predicate classes of interest as

$$C = \{c_1, c_2, \dots, c_{N_c}\}, \quad (1)$$

$$P = \{p_0, p_1, p_2, \dots, p_{N_p-1}\}, \quad (2)$$

where N_c and N_p are the numbers of component and predicate classes in the graph language, respectively. The semantic multi-graph description generated by the detected instances of C and P is defined as

$$\mathcal{D}_t = \{D_{t,1}, D_{t,2}, \dots, D_{t,N}\}, \quad (3)$$

where each $D_{t,i}$ is an atomic bipartite graph that includes a connected set of detected component and predicate instances, as illustrated in Fig. 2. Note that the instance graph includes multiple instances of the same component or predicates; hence, each node in Fig. 2 is identified with a two-tuple, i.e., the component or predicate, and a unique ID number.

Each component and predicate in an atomic graph $D_{t,i}$ has a corresponding attribute set denoted by

$$\Theta_{t,i}(n_j, k) = \{\theta_{t,i}^{(1)}(n_j, k), \theta_{t,i}^{(2)}(n_j, k), \dots, \theta_{t,i}^{(L_{n_j})}(n_j, k)\}, \quad (4)$$

where (n_j, k) is the corresponding component or predicate node in the graph. Note that for each type of node n_j , we define L_{n_j} -levels of attributes. This enables an organized way of storing different attributes, preferably in increasing order of complexity for ease of goal-based filtering. We note that in the computer vision case study given in [9], there are three levels of attributes: the scalar position and velocity in the lowest-level attribute set $\theta_{t,i}^{(1)}$, subfeature vectors in $\theta_{t,i}^{(2)}$, and the full bounding-box images of the detected components in $\theta_{t,i}^{(3)}$. Nonetheless, the definition given in (4) is purposefully general, and can be modified depending on the specific application.

Recent advances in machine learning and artificial intelligence have enabled the extraction of semantic information in the proposed graph structure of (1)–(4). In practice, machine learning and artificial intelligence based semantic extractors do not have perfect accuracy or sensitivity; hence, the extracted semantic information does not have perfect correspondence with the objective reality. Moreover, the evolution of the generated semantic output must be tracked accurately to filter out any erroneous detection and identify the points of innovation so that the transmission (or storage) events can be scheduled efficiently. An advantage of using learning-based algorithms is that statistical properties of the algorithms (e.g., confusion matrices) can be modeled through extensive training statistics, which can be used to improve the fidelity of the semantic output by exploiting the natural temporal continuity of the objects in the sensed area. Additionally, logical probabilities of the possible semantic output patterns can be exploited to improve reliability in a Bayesian sense, as well. For example, in a computer vision security application involving a car dealer

showroom for a certain brand, identifying cars with other brands are much less likely logically, and thus, semantically. Therefore, the semantic extractor can take this into account during the temporal evolution of its semantic output to determine whether any new detections are the result of a true innovation or a sporadic erroneous event.

In the next section, with the objective of addressing the problems stated above, we introduce a semantic extraction framework that is based on the semantic structure detailed in this section. We note that the proposed methods are not limited to the structure described here, and they can also be used for different graph-based signals with embedded scalar/vector attributes.

IV. THE PROPOSED SEMANTIC EXTRACTION FRAMEWORK

In the semantic signal processing framework of [9] summarized in Section III and illustrated in Fig. 1, the semantic extractor is assumed to be generating reliable and smooth semantic output signals. On the other hand, the extracted semantic information includes semantic noise that is introduced by input noise, algorithm limitations, etc. In this section, we propose several methods that can be used individually or collectively to exploit prior information about the semantic extractor and environmental characteristics to improve the fidelity of the extracted semantic information. Moreover, the proposed techniques enable asynchronous updates of the semantic signal across time, where the updates can be designed to occur only when there is a significant semantic innovation in the signal.

The proposed conceptual block diagram for a semantic extractor is given in Fig. 3. Note that the *Raw Semantic Extraction* block typically involves a learning-based algorithm such as a scene-graph generator for image/video signals [53], segmented and classified objects in LIDAR images [79], sound event detection (SED) for acoustic signals [3], etc. The following blocks in Fig. 3 exploit the known characteristics of the raw extractor, such as its confusion matrix; and the statistics of the environment, such as the occurrence probabilities of the possible output signals. Throughout this paper, we assume either these parameters are known a priori or estimated through measurements.

The proposed *Semantic Fidelity Control* block takes the raw output of the initial extraction and uses time integration techniques to improve the detection characteristics according to internally set or externally requested reliability parameters. The proposed fidelity control method is explained in Section V.

Note that after this point in the semantic extraction framework shown in Fig. 3, tracking and smoothing blocks work in parallel processing banks, with each instance working on an atomic graph as defined in (3). This parallel approach is the result of intentional separation of the inter-connected atomic graphs, and reduces the computational complexity of the following blocks. Also note that a list of goals denoted as \mathcal{G}_t is an input of the semantic extraction framework. Either internally or externally defined, the goals can be used to filter the total semantic description of the signal into a goal-oriented

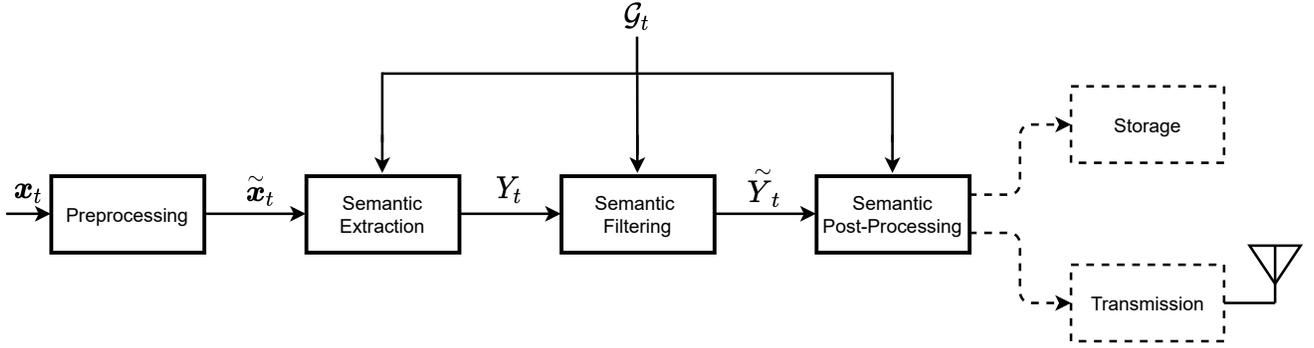


Fig. 1: The proposed goal-oriented semantic signal processing framework [9]. Raw signal input is denoted as x_t , whereas the external goals are denoted with G_t .

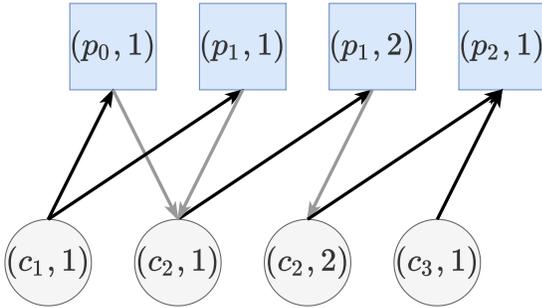


Fig. 2: An example of an atomic bipartite graph structure ($D_{t,i}$) with signal components and predicates at the instance level. The two-tuple representation denotes the component or predicate class, and the unique instance identifier.

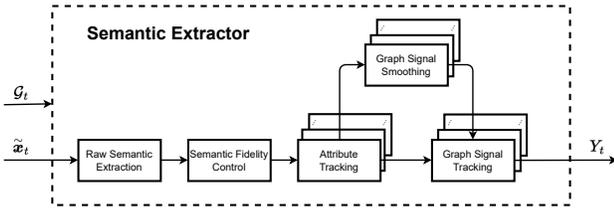


Fig. 3: Proposed block diagram for a typical graph-based semantic extractor.

subset, which can also greatly reduce the processing load of the following blocks. Since the goal-filtering can be performed anywhere along the processing chain, it is left as a general input in Fig. 3.

In a typical application of semantic signal processing or communications, the innovation in the semantic signals will not only come from the graph structure, but also from the low-level numerical attributes embedded in the graphs (see (4) in Section III). The *Attribute Tracking* is an application-specific block, where these low-level numerical attributes can be tracked across time to detect innovations in the semantic signals of interest. For vector attributes (e.g., subfeature vectors or state vectors), we propose and demonstrate a subspace tracking algorithm in Section VI.

The *Graph Signal Smoothing* block is proposed as an optional pre-filtering step to smooth out sporadic erroneous

detections and reduce the throughput for the following blocks. Hence, this step needs to have a relatively low computational cost compared to the final filtering/tracking block. Therefore, a simple GED algorithm with a modified cost definition is proposed for this block and explained in Section VIII.

The final block in Fig. 3, the *Graph Signal Tracking*, uses the pre-filtered graph signal and attribute updates in conjunction with the estimated environmental probabilities to produce a reliable, accurate semantic description of the input signal. For this block, we propose an HMM over the graph signals in Section VII and use state estimation algorithms to optimize the output signal. As a result, the output of this block is a reliable, accurate, and slowly-varying signal (compared to the raw input) that can indicate significant innovations in the graph-signal or its underlying attributes.

Starting with the next section, we present each sub-block and the potential methods and algorithms that can be incorporated within, in detail. Note that the proposed algorithms and techniques presented in the following sections can be used collectively, as shown in Fig. 3, or individually depending on the computational capabilities and requirements of the signal processing hardware.

V. SEMANTIC FIDELITY CONTROL WITH TIME INTEGRATION

Time integration of signals for improved detection and estimation performance are well documented in the signal processing literature [80]–[82]. These techniques can also be employed in the semantic signal processing framework in the initial *preprocessing* block, shown in Fig. 1, in front of the *semantic extractor*. With the introduction of a semantic language and a semantic extractor, a semantically-aware way of improving the detection and estimation performance can be achieved by employing time integration techniques at the output of the semantic extractor.

The *raw semantic extractor* shown in Fig. 3 typically includes a machine learning algorithm that detects the signal components and their inter-relationships. After the training and testing of these algorithms, we assume that we have access to the confusion matrix that provides statistics on the

characteristics of the trained algorithm [83]. The confusion matrix of a raw extractor is in the following form:

$$CM(\tau) = \begin{bmatrix} n_{1,1} & n_{1,2} & \cdots & n_{1,K} \\ n_{2,1} & n_{2,2} & \cdots & n_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ n_{K,1} & n_{K,2} & \cdots & n_{K,K} \end{bmatrix}, \quad (5)$$

where τ is the score threshold used for detection, $n_{i,j}$ corresponds to the number of times a pattern j is detected given an actual pattern of i , with K being the total number of output patterns. The output patterns in this definition can be simple components/predicates or more complex graph outputs as shown in Fig. 2, i.e., the range of possible semantic outputs that the extractor can generate. Note that, the total number of possible components and predicates (N_c and N_p) are limited, since they are defined for a specific application with specific goals. Moreover, in a given time instant, the connected atomic semantic graph outputs will have even fewer components and predicates (η_c and η_p , respectively), which will reduce the computational complexity of processing the graphs even further.

The confusion matrix can be used to infer other statistical metrics about the output patterns. Let $N_i \geq \sum_j n_{i,j}$ be the total number of samples where the pattern i exists, and $N = \sum N_i$ be the total number of samples. With (5) and N_i , other pertinent confusion metrics for each pattern can be written as

$$\text{True Positive Rate: } TPR_i = \frac{n_{i,i}}{N_i}, \quad (6)$$

$$\text{False Positive Rate: } FPR_i = \frac{\sum_{j \neq i} n_{j,i}}{N - N_i}, \quad (7)$$

$$\text{False Negative Rate: } FNR_i = \frac{N_i - n_{i,i}}{N_i}, \quad (8)$$

$$\text{True Negative Rate: } TNR_i = \frac{\sum_{j \neq i} \sum_{k \neq i} n_{j,k}}{N - N_i}. \quad (9)$$

Note that (6)–(9) also depend on the detection threshold τ . By changing τ , different detection characteristics for a given pattern can be obtained to generate the Receiver Operating Characteristic (ROC) curve. The ROC for the semantic extractor can also be constructed for different environmental factors such as time of day, weather, season, etc., depending on the application and the desired accuracy.

To control the fidelity of the extractor output with a semantically-aware approach, we propose temporal integration of the output scores of a classifier that identifies the signal components to improve the ROC. Assuming underlying Gaussian processes, the output score distribution (S_i) for each pattern under binary hypothesis testing (\mathcal{H}_0 : the pattern does not exist and \mathcal{H}_1 : the pattern exists) can be modeled as

$$P(O_i) \sim \begin{cases} \mathcal{N}(\mu_{i,0}, \sigma_{i,0}^2), & \text{under } \mathcal{H}_0 \\ \mathcal{N}(\mu_{i,1}, \sigma_{i,1}^2), & \text{under } \mathcal{H}_1 \end{cases} \quad (10)$$

where μ and σ are the mean and standard deviation for the score distributions in each case. Assuming a jointly Gaussian

distribution, the time-average of T_i consecutive samples yields

$$P(O_i^{T_i}) = P\left(\frac{1}{T_i} \sum_{n=1}^{T_i} O_i(t_n)\right) \sim \begin{cases} \mathcal{N}(\mu_{i,0}, \hat{\sigma}_{i,0}^2), & \text{under } \mathcal{H}_0 \\ \mathcal{N}(\mu_{i,1}, \hat{\sigma}_{i,1}^2), & \text{under } \mathcal{H}_1 \end{cases} \quad (11)$$

where $O_i(t_n)$ is the output score at time t_n . As expected, the mean values do not change and the standard deviations are bounded as

$$\sigma_i/T_i \leq \hat{\sigma}_i \leq \sigma_i, \quad (12)$$

for both hypotheses. The lower and upper bound for the standard deviations are achieved when the consecutive samples are independent or perfectly correlated, respectively. To illustrate the correlation between consecutive samples, consider the semantic extraction of a video signal. If the scene is perfectly stationary, i.e., the same still images are recorded across consecutive time instants, the score distributions will be identical, rendering the time integration ineffective. However, even small changes in the video image can lead to different score output realizations of the same underlying distribution, hence improving the detection performance.

Since the standard deviations for the distributions under both hypotheses are reduced, it is straightforward to show that the ROC metrics given in (6)–(9) can be improved as

$$TPR_{i,T_i} \geq TPR_i, \quad (13)$$

$$FPR_{i,T_i} \leq FPR_i, \quad (14)$$

$$FNR_{i,T_i} \leq FNR_i, \quad (15)$$

$$TNR_{i,T_i} \geq TNR_i, \quad (16)$$

with subscript T_i denoting the time-integrated metrics. Again, the equalities in (13)–(16) hold only when the consecutive samples are perfectly correlated. Also, note that the assumption of Gaussianity is not essential, i.e., similar expressions can be written for other distribution models as well.

To illustrate the effects of integrating the output scores across time, we have built a test dataset based on the images from the *Cars Dataset* given in [84], where we placed scaled versions of the car images randomly onto a stationary parking lot background image, as shown in Fig. 4. A total of 897 cars randomly selected from the *Cars Dataset* were placed randomly onto the background shown in Fig. 4 at different size scales, each for a total of 100 realizations. Then, the YOLOv4 algorithm [20] is used to generate the output scores for each case, resulting in the output scores and the ROC curves shown in Fig. 5. As illustrated in Fig. 5, the output scores and the corresponding detection performance deteriorate monotonically with the decreasing size of the object. Using the output scores generated by YOLOv4, we integrate different realizations of the same car objects for various integration window lengths. The resulting improvement in the ROC curves for scales 5% and 15% are illustrated in Figs. 6 and 7, respectively, which clearly show the expected monotonic increase in the detection performance as the integration window length increases.

Integration across time also means a delay in producing the desired detections. Hence, depending on the desired FPR and maximum allowable *time-on-target*, the raw score output can be integrated to achieve the maximum possible TPR . With



Fig. 4: An example car image placed onto a stationary parking lot background. The scale number annotation corresponds to the ratio of the size of the car compared to the full image.

a semantically-aware approach, these parameters can be independently defined or estimated for different types of targets (e.g., time-on-target can be a function of the coherence time and the maximum delay) and environmental factors as well. This idea will be illustrated via case studies in Section IX.

VI. SUBSPACE TRACKING OF VECTOR ATTRIBUTES

The overall goal of the semantic extractor framework shown in Fig. 3 is to provide reliable and smooth semantic signals while also identifying time indices of significant innovation for scheduling transmission or storage events. The potential fidelity improvement provided by the time integration method only works over short time periods where the semantic information content within the signal does not change (analogous to coherence time in wireless channels) significantly. The following blocks in Fig. 3, starting with the *attribute tracking* aim to smooth the semantic signal and identify points of significant innovation over extended periods of time.

In the raw semantic extraction block shown in Fig. 3, the semantic data are organized to include goal-filtered components and their corresponding attributes. In a typical computer vision problem, these attributes can be position, speed for scalar attributes, sub-feature vectors for vector attributes, or even the encoded video signal itself. Tracking these attributes enables identification of strong changes at the attribute level, which in turn can be used to update the transmitted or stored attributes, even when the overarching semantic graph structure stays the same. On the other hand, if the semantic graph signal is erroneously updated (or kept the same), attribute tracking may be used to reconcile consecutive attributes to detect and correct these errors.

We now present a subspace tracking algorithm that can be used for vector attributes. To illustrate the algorithm more clearly, we use the real-world semantic extraction of video-streams presented in [9] where we define signal components as detected objects in the stream. As illustrated in Fig. 8, the proposed semantic extractor in [9] utilizes YOLOv4-CSP [20] model to detect occurring objects in each input frame F_n . For temporal extension, we adopt *tracking by detection* paradigm

and utilize DeepSORT [85] to track individual objects in the video stream which utilizes Kalman filter to recursively predict future positions of detected objects. Furthermore, DeepSORT inherits another small convolutional neural network model trained on re-identification task to extract visual feature vectors of each detected object to be used for association of existing tracks and recent detections. We refer readers to [85]–[88] for details about tracking by detection paradigm. We denote these feature vectors with $r_i \in \mathbb{R}^{128}$, where $\|r_i\|_2 = 1$, and utilize them as second-level vector attributes of component nodes in our semantic graph output (see [9] for details). We would like to point out that the feature vector r_i for a particular detected component is a 128-dimensional vector and carries most of the attribute-level *innovation* in the semantic signal processing framework. In Fig. 8, the object-level goal \mathcal{G}_{object} restricts the components and the predicates that must be detected. The corresponding *multi-graph instance representation* $\mathcal{D}_n = \{D_1, \dots, D_{M_n}\}$ composed of M_n disconnected subgraphs and a corresponding attribute set $\mathcal{A}_n = \{A_1, \dots, A_{M_n}\}$ extracted from the video signal. The *multi-graph class representation* \mathcal{S}_n is constructed in the Graph Abstraction block, and the complete semantic description $Y_n = (\mathcal{S}_n, \mathcal{D}_n, \mathcal{A}_n)$ is generated. Finally, if there are external goals \mathcal{G}_n , the semantic description is filtered to obtain the goal-filtered semantic description.

To process and track attribute vectors efficiently, we use dimensionality reduction techniques such as Robust PCA and subspace tracking [89]. To track the evolution of the subspace of feature vectors with respect to time, we use Fast Principal Component Pursuit (PCP) via alternating minimization [90].

Defining a feature matrix D , with its rows being the subfeature vectors across time, and decomposing the matrix into a compressed basis matrix L and residual matrix S , the PCP problem is formulated as

$$\arg \min_{L,S} \|L\|_* + \lambda \|S\|_1 \quad \text{subject to} \quad D = L + S \quad (17)$$

where $\|L\|_*$ is the nuclear norm of matrix L and $\|S\|_1$ is the l_1 -norm of matrix S . A variant of (17) can be constructed as

$$\arg \min_{L,S} \frac{1}{2} \|L + S - D\|_F + \lambda \|S\|_1 \quad \text{subject to} \quad \text{rank}(L) = t. \quad (18)$$

Finally, we can solve the problem via alternating minimization, namely,

$$L_{k+1} = \arg \min_L \|L + S_k - D\|_F \quad \text{subject to} \quad \text{rank}(L) = t \quad (19)$$

$$S_{k+1} = \arg \min_S \|L_{k+1} + S - D\|_F + \lambda \|S\|_1. \quad (20)$$

The only computationally demanding step is (19), and it can be solved by computing a partial SVD of $D - S_k$ with t components. The solution to (20) is *element-wise shrinkage*, i.e., $\text{sign}(D - L_{k+1}) (|D - L_{k+1}| - \lambda)^+$, where $(x)^+ = \max\{0, x\}$. Note that t in (19) can be chosen by a heuristic procedure. We increment t and estimate the contribution of the new singular vector. Comparing the contribution with a threshold, the algorithm stops at a particular t . In numerical simulations, we observed that t is usually small (does not exceed 3 in our setting).

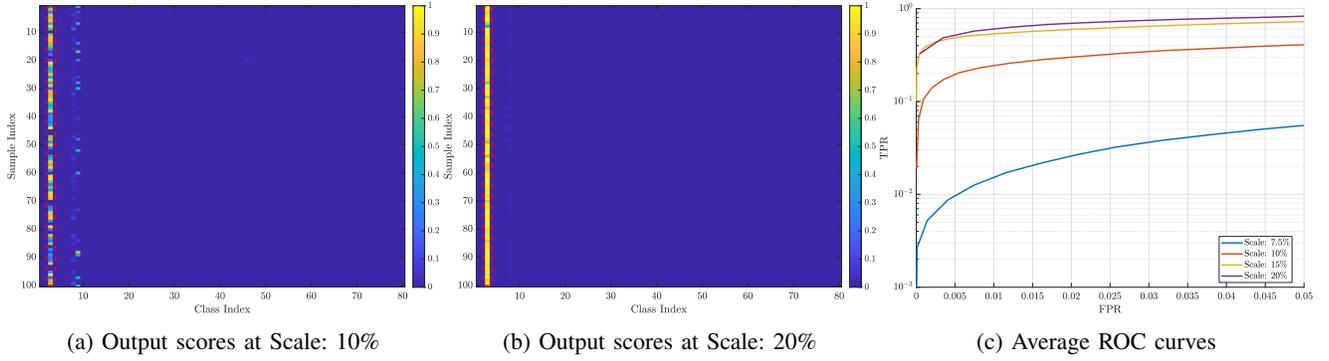


Fig. 5: Output scores and the ROC curves generated by YOLOv4 [20]. Note that the correct class detection for the *car* class is highlighted in (a) and (b).

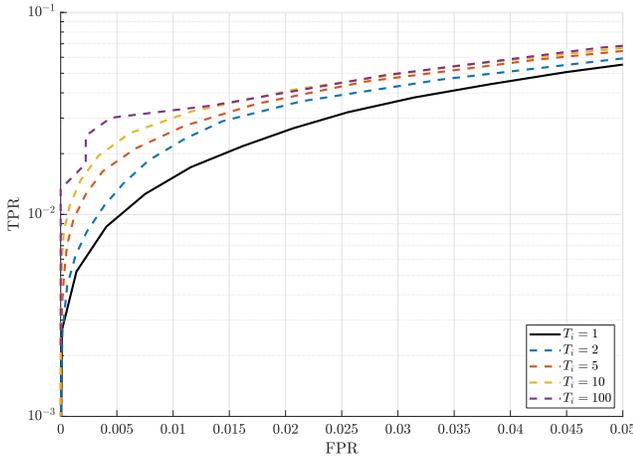


Fig. 6: ROC curves for 5% scaled car images at different integration lengths.

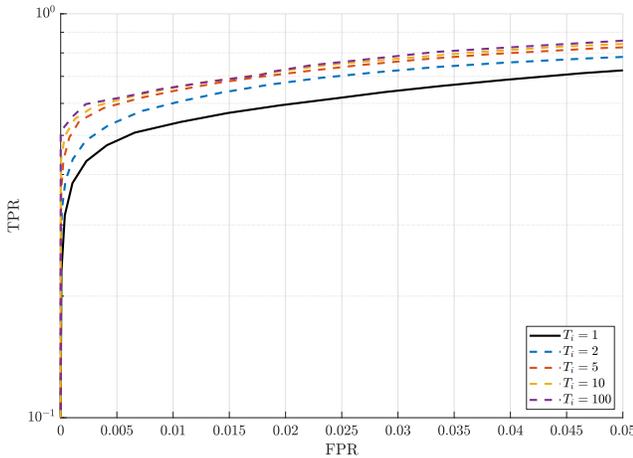


Fig. 7: ROC curves for 15% scaled car images at different integration lengths.

To utilize PCP to track the evolution of the feature vectors in a practical application, one can solve the alternating minimization problem given in (19) and (20) periodically on a moving window. The size of this window (buffer) depends on the application and the rate of change in the semantic content.

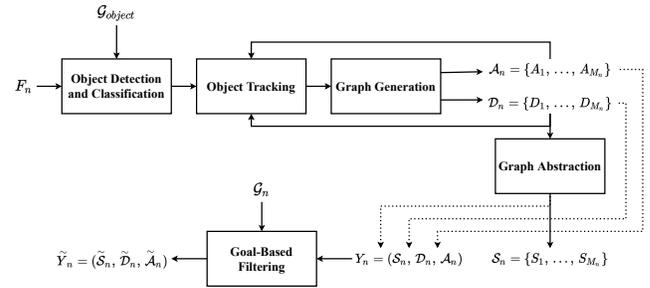


Fig. 8: Semantic extraction for real-time video signals as presented in [9].

Then, a simple heuristic approach can be taken to detect innovations by checking the l_1 -norm of the sparse component S after convergence and compare it with a predefined threshold.

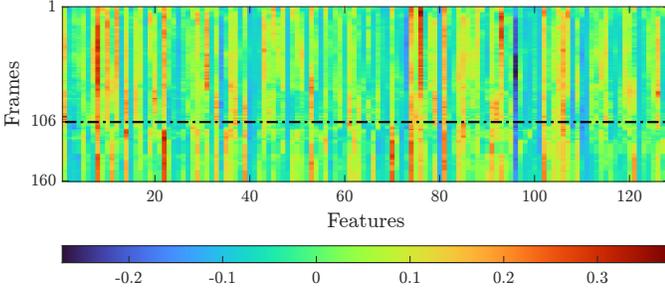
To illustrate the inherent redundancy for consecutive frames, we show the still images from a 10s video that includes a car passing from a shaded region to a sunny region in Fig. 9. The correlations of consecutive attribute vectors are given in Fig. 10a. Note that in the scenario given in Fig. 9, the semantic components and predicates stay the same (i.e., the same objects and relationships are present throughout the clip), while the attributes of the car component (car-3) change due to the image brightness, contrast, etc. As shown in Fig. 10a, an innovation at the attribute level can be identified around frame-106, when the lighting changes for the car-3 component in the overall semantic graph signal. In Fig. 10b, the attribute level innovation has been illustrated by analyzing the l_1 -norm of the sparse component of the feature matrix D . As the norm is maximum at frame 106, we identify the innovation and locate the precise time instance quantitatively.

VII. GRAPH SIGNAL TRACKING USING A HIDDEN MARKOV MODEL

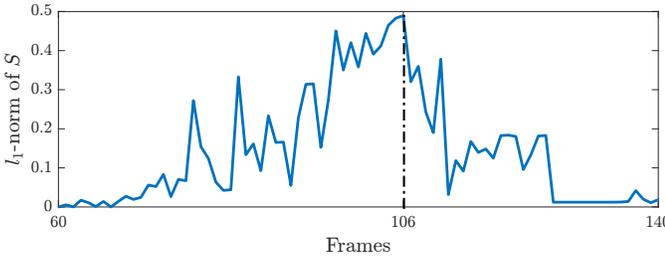
After the *semantic fidelity control* and *attribute tracking* blocks, we now have a more reliable signal in which we can detect significant innovations at the attribute level. To provide a smooth output while also identifying significant innovations at the graph level, we introduce a stochastic model for the temporal evolution of the graph signals. More



Fig. 9: A car is passing from a shaded region to a sunny region. The middle and right sub-figures demonstrate the *innovation* event at the attribute level, while the semantic components stay the same.



(a) Correlation of subfeature vectors across consecutive frames. The innovation at the attribute level is illustrated by a black dashed-dotted line.



(b) Frames vs. l_1 -norm of the sparse component S . Note that the attribute level innovation is visible around the frame-106 since the l_1 -norm of the sparse component achieves its maximum.

Fig. 10: Attribute level innovation analysis for the car example given in Fig. 9

specifically, a Markov Chain can be built by defining distinct component-predicate connections in the graph representation in Section III as its distinct states. However, these states are not directly observable, as there may be erroneous detections due to semantic or technical noise. Therefore, we augment the Markov chain with an HMM, where the observed output sequence is generated by the semantic extractor as a function of the actual states of the underlying Markov chain.

The HMM assumption of the graph representation enables estimation of the underlying state sequence through the output sequence generated by the previous blocks. This stochastic model also enables quantifying the semantic information through the entropy rate of the Markov chain, and it can be used to compress the semantic signals. Throughout this section, we rigorously define the HMM and explain how

the underlying state sequence can be extracted from the observations.

A. Markov Chains and the Hidden Markov Model

As shown in Fig. 2, a graph signal can have various connections among component and predicate nodes. Each possible graph configuration can be defined as distinct states of the graph. Since the number of component and predicate nodes in the graph is finite (and very small in the proposed language structure of [9] with goal filtering), the number of possible graph states is limited. The state-space of the graph can be defined as $S = \{S_0, S_1, \dots, S_{N-1}\}$. With the Markovity assumption, the state transition probability of the graph can be expressed as P_{S_j, S_i} for making a transition from S_j to S_i in a single frame

$$P_{S_j, S_i} = P(Q_{t+1} = S_i | Q_t = S_j, Q_{t-1} = S_{Q_{t-1}}, \dots, Q_0 = S_{Q_0}) = P(Q_{t+1} = S_i | Q_t = S_j), \quad (21)$$

where Q_t represents the state variable at time instance t .

The temporal evolution of the graph can be characterized entirely through its state transition matrix (\mathbf{A}) and the initial probability distribution (\mathbf{p}). \mathbf{A} is defined such that its (j, i) 'th entry is equal to P_{S_j, S_i} , which is the probability of moving to state S_j from state S_i in a single step. The i 'th element of (\mathbf{p}), denoted by p_i , which is the probability of being in state S_i in the beginning.

The HMM of the graph consists of two sequences: state sequence (\mathbf{Q}) and observation sequence (\mathbf{O}). The state sequence (\mathbf{Q}) consists of actual states, and the output sequence (\mathbf{O}) is the observed state sequence. Elements of both sequences are from the state set S . At each time instance, the observed state can be the actual state, or it can be replaced with another element from the state set, leading to an incorrect observation.

In an HMM, the graph evolves as a Markov Chain according to its state transition matrix (\mathbf{A}) and initial state distribution (\mathbf{p}). As the graph moves between different states (Q_t) over time, the observed state (O_t) changes accordingly. O_t is a random variable which is a function of the actual state (Q_t) at time t . Probability of observing S_j when the underlying state is S_i at time t is denoted by P_{S_i, O_j} , i.e.,

$$P_{S_i, O_j} = P(O_t = S_j | Q_t = S_i). \quad (22)$$

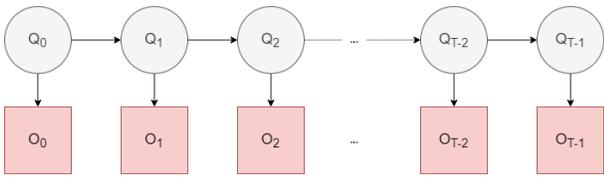


Fig. 11: Hidden Markov Model: Q_t is the actual state and O_t represents the observed state at time t . O_t depends solely on Q_t .

The HMM corresponding to the underlying semantic information can be represented completely through the state transition matrix (\mathbf{A}), initial probability distribution (\mathbf{p}) and the observation matrix (\mathbf{B}). The (i, j) 'th element of (\mathbf{B}) is equal to P_{S_i, O_j} , which is the probability of observing S_j when the actual state is S_i .

As previously mentioned, a major motivation behind the stochastic model is to solve the smoothing problem by estimating the true nature of the observed process through the noisy output of the semantic extractor. To solve the smoothing problem, we are interested in estimation of the underlying state sequence (\mathbf{Q}) from the observed state sequence (\mathbf{O}) when the model is completely known. This is a well-known problem in HMM framework [91]; however, the definition of optimality for the solution varies. There are two common approaches: The first one is to find the most likely state sequence producing the observed state sequence (\mathbf{O}), for which the state sequence (\mathbf{Q}) can be estimated using the Viterbi Algorithm [92]. The second approach is to find the state sequence consisting of the most likely states at each step, maximizing the number of correctly estimated states on average. The solution to the second formulation follows from the Forward-Backward Algorithm [93].

The goal of the proposed framework is to find the most likely state sequence producing the observed state sequence; hence, we focus on the application of the Viterbi Algorithm. Viterbi Algorithm is a Dynamic Programming based recursive algorithm that finds the state sequence maximizing the posterior probability given the output sequence and model parameters ($P(Q|O; \mathbf{A}, \mathbf{B}, \mathbf{p})$).

Some definitions are required to describe the algorithm fully. The underlying state sequence is represented as $Q = \{Q_0, Q_1, \dots, Q_{T-1}\}$ where T is the sequence length. The observed state sequence corresponding to (Q) is denoted as $O = \{O_0, O_1, \dots, O_{T-1}\}$. The model parameters: state transition matrix (\mathbf{A}), observation matrix (\mathbf{B}), and initial state distribution (\mathbf{p}) are represented compactly as λ . At each time step t , Viterbi Algorithm considers every state i visited. Then, the probability of observing $\{O_0, O_1, \dots, O_t\}$ and terminating the path at state i $\{Q_t = S_i\}$ is maximized over the previous states $\{Q_0, Q_1, \dots, Q_{t-1}\}$ [92], i.e.,

$$\delta_t(i) = \max_{\{Q_0, Q_1, \dots, Q_{t-1}\}} P(Q_0, Q_1, \dots, Q_{t-1}, Q_t = S_i, O_0, O_1, \dots, O_t | \lambda). \quad (23)$$

The heart of the Viterbi Algorithm is the induction step

which is used to find $\delta_{t+1}(j)$ from $\delta_t(i)$ as

$$\delta_{t+1}(j) = [\max_{S_i} \delta_t(i) a_{i,j}] b_{j, O_{t+1}}, \quad (24)$$

where $\delta_t(i)$ denotes the probability of most likely state sequence $\{Q_0, Q_1, \dots, Q_{t-1}\}$ ending at $Q_t = S_i$ given the observation sequence $\{O_0, O_1, \dots, O_{t-1}\}$ for the first t steps. To find the most likely state sequence, each state for each time instance must keep a pointer $\psi_t(j)$ to the previous state maximizing (24). The most likely state sequence is found by backtracking the pointer of the state maximizing $\delta_{T-1}(i)$. The complete procedure is summarized in Algorithm 1.

Algorithm 1 Viterbi Algorithm

Input: $\mathbf{A}, \mathbf{B}, \mathbf{p}, \mathbf{O}$ \triangleright Model parameters and output sequence
Output: \mathbf{Q} \triangleright Most likely state sequence

Initialization
for $i = 0$ upto $N-1$ **do**
 $\delta_0(i) = p_i b_{i, O_0}$
 $\psi_0(i) = -1$
end for

Induction
for $t = 1$ up to $T-1$ **do**
for $j = 0$ up to $N-1$ **do**
 $\delta_t(j) = \max_i [\delta_{t-1}(i) a_{i,j}] b_{j, O_t}$
 $\psi_t(j) = \arg \max_i [\delta_{t-1}(i) a_{i,j}]$
end for
end for

Termination
 $P^* = \max_i [\delta_{T-1}(i)]$ \triangleright Probability of most likely state sequence
 $Q_{T-1}^* = \arg \max_i [\delta_{T-1}(i)]$

Backtracking
for $t = T-2$ down to 0 **do**
 $Q_t^* = \psi_{t+1}(Q_{t+1}^*)$ \triangleright Backtrack the most likely path
end for

return Q^* \triangleright Return the most likely state sequence

The Viterbi Algorithm finds the most likely state sequence from the observations in $\mathcal{O}(N^2T)$ steps, where N denotes the number of states and T is the sequence length. Hence, the computation of the algorithm can become intractable for very large number of states. However, in our proposed hierarchical graph representation explained in Section III, we separate the semantic graph description of the raw signal into connected atomic graphs that are expected to have a very limited number of states (components and predicates), especially considering that the proposed semantic graph description already starts with a limited component and predicate set for a specific application. Another dramatic reduction over the complexity of a practical implementation is introduced by the goal-oriented nature of the proposed framework, i.e., only a subset of the whole semantic description of a signal is of interest, and is being processed, stored, or transmitted, as discussed in Section IV.

Even further reductions in the complexity of an HMM implementation is achievable with the use of approximate algorithms such as the M-Algorithm, which is a greedy

algorithm that considers only the most likely M states as the predecessor during the induction step. M-Algorithm has a complexity of $\mathcal{O}(MNT)$ where M is strictly less than N [94].

Aside from the computational complexity, another important aspect of successfully implementing an HMM-based algorithm is the accurate estimation of the model parameters. In this work, we assume that the parameters of HMM such as the state transition probabilities are known a priori. For practical applications, the model parameters can be initialized with available logical prior information (some state descriptions can be logically impossible etc.), then can be estimated from the observations to provide an accurate model for HMM. The Baum-Welch algorithm is a popular method for the model estimation problem in the HMM implementations [95]. The algorithm starts with an initial guess on the model parameters. It finds the new parameters iteratively by maximizing the observation likelihood for the estimated parameters at the current step. This two-step procedure is repeated until the model parameters converge.

Another potential shortcoming of the proposed HMM setting is the restrictions on the waiting time distributions. A Markov model enforces each state waiting time to follow a geometric distribution. In many practical cases, the observed semantic signal might contain states that cannot be modeled accurately with a geometric distribution, leading to a discrepancy between the underlying semantic signal and its Markov model. To avoid such discrepancies, HMM can be replaced with a Hidden Semi-Markov Model (HSMM). Unlike HMM, a state in HSMM might have any distribution for its waiting time, enabling the modeling of a broader range of semantic signals in practice, albeit with a higher computational cost. The observed graph sequence, which follows a HSMM, can be smoothed with a slightly modified version of the Viterbi Algorithm [96].

VIII. GRAPH SIGNAL SMOOTHING USING GED

In the proposed semantic extraction framework, the *Graph Signal Tracking* block uses an HMM-based estimation algorithm to smooth and track the semantic graph signals. However, due to the computational complexity of the state sequence estimation procedure given in Algorithm 1, a relatively simple algorithm to identify semantic innovation and smooth erroneous detections is introduced in this section as part of the *Graph Signal Smoothing* block in Fig. 1. The algorithm and the smoothing block presented here can either be used as a pre-filtering step before the *Graph Signal Tracking* block to reduce the computational load on the HMM state estimation algorithm, or it can be used as an alternative to the HMM for a low complexity implementation.

To identify actual or erroneous changes in the graph signals with relative simplicity, we use a graph distance metric for consecutive graphs across time. GED [97], [98] is one such popular metric. For two graphs g_1 and g_2 , the GED is defined as follows:

$$GED(g_1, g_2) = \min_{(e_1, e_2, \dots, e_k) \in \mathcal{P}(g_1, g_2)} \sum_{i=1}^k c(e_i), \quad (25)$$

where $\mathcal{P}(g_1, g_2)$ is the set of all edit paths that transform g_1 into g_2 , e_i 's are the elementary edit operations, and $c(e_i)$ is the cost of performing the edit e_i . Elementary edit operations include the insertion, deletion, and substitution of nodes and edges. Note that with the semantic language defined in Section III, we have a bipartite graph structure whose edges are always strictly dependent on the predicates that connect different components. Therefore, for the proposed language, we need to define edit costs only for the nodes of the bipartite graphs.

In most applications of the GED, the cost functions are chosen to be constant scalars depending on the type of operation [98], [99]. In this work, we define the edit costs to be a function of the confusion metrics of the underlying semantic extractor. More specifically, we exploit the known statistical characteristics of the semantic extractor that generates the graphs to obtain the GED as a statistical similarity between two consecutive graphs. As such, we define the cost as

$$c(e_i) = -\log(P(e_i)), \quad (26)$$

where $P(e_i)$ is the probability of the edit e_i occurring at the output of the semantic extractor. Therefore, the GED becomes an estimator for the probability of obtaining the new graph g_2 , given the initial graph g_1 . Note that this probability can be estimated using the confusion matrix and statistical metrics given in (5)–(9). The relationships between elementary edits and corresponding confusion metrics are given in Table I. Note that, if the occurrence probabilities of certain graph realizations are available, posterior probabilities can be used as the distance metrics using Bayesian inference.

TABLE I: GED Cost Definitions

Elementary Edits	Confusion Metric
Node Insertion - i	FPR_i
Node Deletion - i	FNR_i
Node Substitution $i \rightarrow j$	$n_{i,j}/N_i$
Node Substitution $j \rightarrow i$	$n_{j,i}/N_j$

If the confusion metrics for the composite edit paths between g_1 and g_2 are available, they can be used to estimate the GED cost directly. However, if only elementary edit confusion metrics are available, the GED given in (25) becomes a rough estimate since this definition assumes that the individual edits that make up an edit path are statistically independent, as the log-probabilities add up to form the total GED. As the numerical examples and the experimental results show, this assumption is still accurate enough to track the semantic information and identify points of significant innovation. In Section IX, we present several case studies involving simulations and experimental results that showcase the potential GED as an elegant pre-filtering step to filter semantic noise and provide indications of strong semantic innovation.

IX. CASE STUDIES

In this section, we present case studies for each of the techniques through simulations and experimental results. For showcasing the methods that rely on the semantic extractor's confusion matrix without having to limit ourselves to

1	6980	182	7	93	676	87.9%	12.1%
2		5858		101	60	97.3%	2.7%
3	29		6359	374	697	85.3%	14.7%
4			1	1688	4	99.7%	0.3%
5	57	20	232	39	2565	88.1%	11.9%

	98.8%	96.7%	96.4%	73.6%	64.1%
	1.2%	3.3%	3.6%	26.4%	35.9%
	1	2	3	4	5

Predicted Class

Fig. 12: Confusion matrix from a randomly generated semantic extractor with $\tau = 0.5$. The additional tables at the right and bottom side denote the accuracy rates of each row and column, respectively.

a potentially limited training set, we have constructed a simulation framework that can generate a semantic extractor model with random characteristics. This extractor is then used to produce noisy semantic signals, given a randomly generated time evolution of a noiseless semantic signal as the ground truth. For experimental demonstrations, we use the semantic extractor given in [9], which is a YOLOv4-based algorithm that generates noisy semantic signals from video signals. The semantic extractor in [9] uses as its component list the YOLOv4 class labels (see (1)). Then, three predicates (*exists*, *moving together*, *conjunct*) are defined and detected among the components to generate the full bipartite graph outputs. Each component node in the semantic graph also has an embedded attribute list with $L_{n_j} = 3$ according to (4). The first level attribute contains the bounding box position and velocities, the second level attribute contains vector subfeatures of length 128, and the third level attribute includes the raw image of the detected bounding box. More details on these definitions can be found in [9].

Using the extractor described above, we generate noisy semantic signals from several short clips that we recorded. These noisy signals are then post-processed with the proposed methods, and their performance is demonstrated.

A. Semantic Fidelity with Time Integration

Using the simulation framework described above, we generate a random semantic extractor that can detect five different semantic patterns. The normalized output scores of this random extractor are thresholded at different levels to generate the confusion matrices and the ROC curves for its output patterns. An example confusion matrix with the threshold set at $\tau = 0.5$ is shown in Fig. 12. With the semantic extractor characterized in Fig. 12, we generate a random time series of inputs of length 10,000 as shown in Fig. 13 and the corresponding output scores. Then, a moving integration window of lengths 1-to-

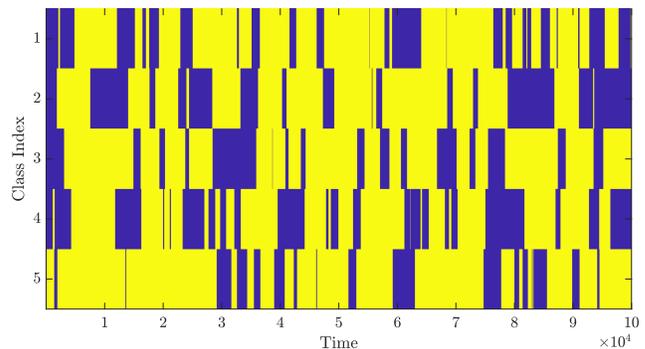


Fig. 13: Randomly generated time evolution for the semantic ground truth. Yellow regions denote where the patterns exist.

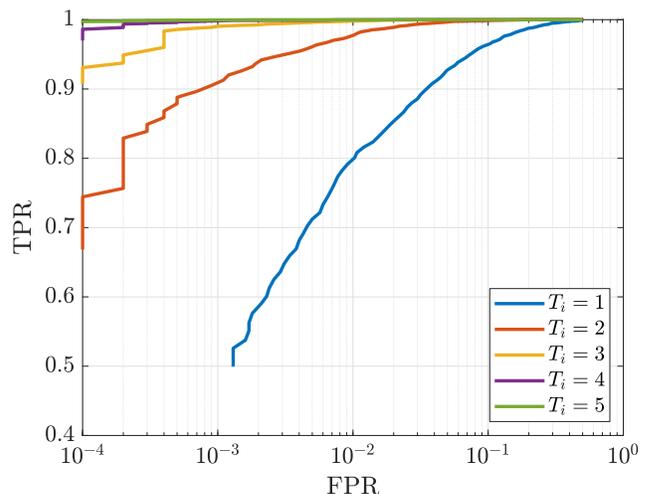


Fig. 14: Improvements in the ROC curves of Class-3 using time integration of the output scores at different window lengths.

5 is implemented to observe the effects of time integration on the detection performance. In Figs. 14 and 15, we present the improvements in ROC curves for the detection of class-3 (C_3). As seen in both graphs, time integration of the extractor output greatly improves the detection performance in this particular scenario. Note that the integration window lengths, TPR, and FPR values presented here are to give a general idea on how time integration affects the performance of the extractor. In practice, these values depend on the input signal acquisition characteristics and the time evolution of the actual semantic content of the signals. On a case-by-case basis, parameter surveys as shown in Fig. 15 can be performed to decide on a *time-on-target* parameter to maximize the detection performance in a class-aware fashion.

B. Attribute Subspace Tracking

To illustrate the vector subspace tracking method developed in Section VI, we use the same recorded video example shown in Fig. 9. Running the PCP algorithm given in (17)–(20) on the video clip, we obtain the subspace breakdown of vector attributes shown in Fig. 16. As it is seen in the sparse component plot, we have two innovations, i.e., when the car

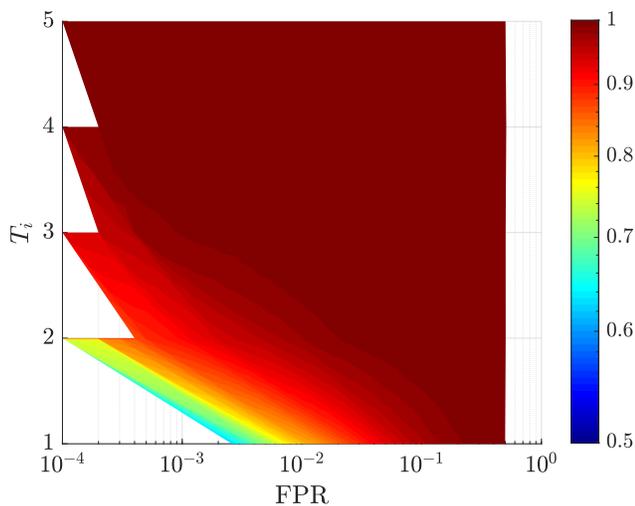


Fig. 15: Heatmap of TPR values for different time integration windows and FPR values.

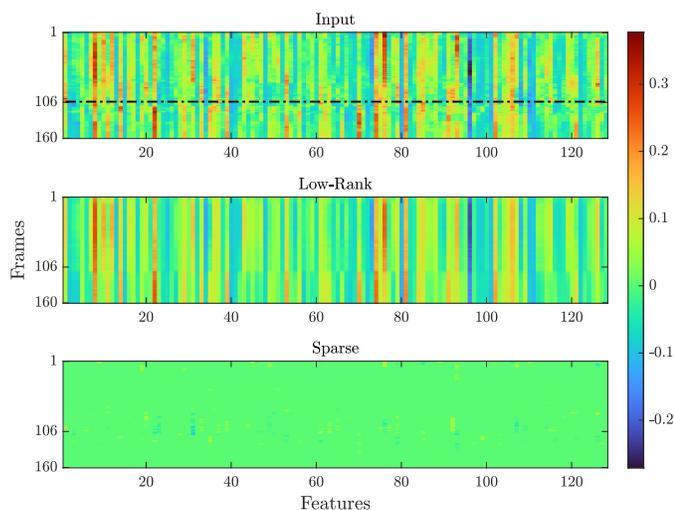


Fig. 16: PCP output of the car video

first appears in the scene and when it passes to the shaded region (around frame 106). By plotting the l_1 -norm of the sparse component, we can track the innovation versus time.

To design an attribute tracking algorithm that can work in real-time, one can use fixed-sized buffers and run the algorithm continuously. Again, we can track the l_1 -norm of the sparse components to detect an attribute-level innovation in the scene, as shown in Fig. 17. Moreover, Fig. 18 illustrates the effect of buffer size on the innovation detection. As expected, smaller buffer sizes result in highly concentrated, spike-shaped innovations. Furthermore, the innovation at the beginning diminishes as the buffer size gets smaller.

In certain scenarios where multiple instances of an object are classified with unique identifiers (ID), the same instance of the class can be detected; however, ID numbers can change due to occlusions, buffer sizes, or missed detections. We demonstrate that even in these cases, the proposed method is able to reconcile the different identifications. In Fig 19, an object detected with multiple IDs is shown throughout the

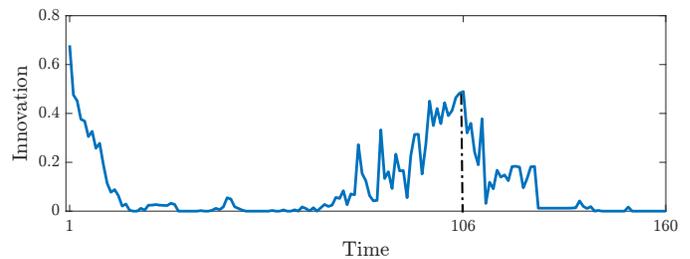


Fig. 17: Innovation vs time for the car video

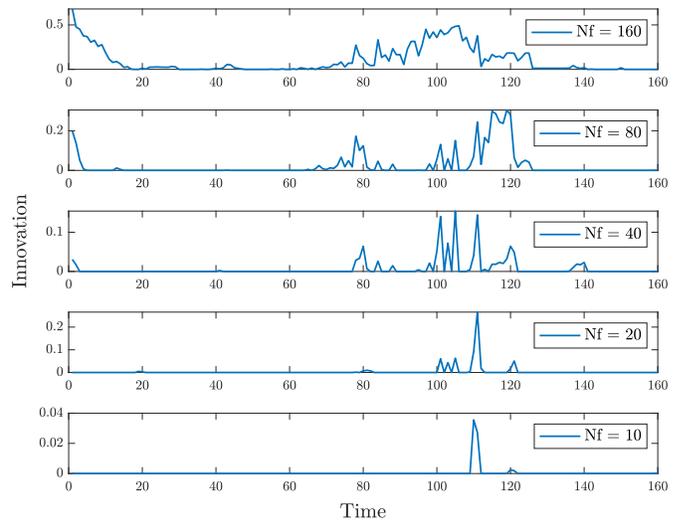


Fig. 18: Innovation vs time for various buffer sizes for the car video

video. Using the PCP algorithm, we analyze the low rank components of these seemingly separate IDs in Fig 20. To reconcile these objects, we calculate the mean of low-rank components over time and compare the Manhattan distance between them, e.g., $\|L_7 - L_1\|_1$ vs. $\|L_7 - L_2\|_1$ to reconcile the instance with ID 7. The resulting distance comparison is given in Fig 21. As illustrated, IDs 7, 16, and 17 are correctly reconciled with the instance ID 2.



Fig. 19: An object is assigned different IDs during the video

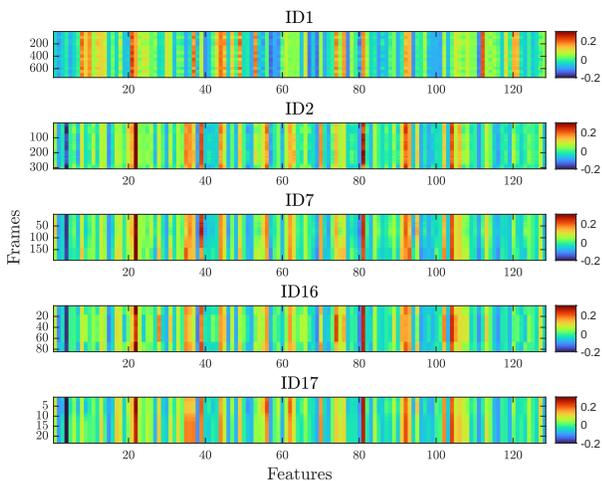


Fig. 20: Low-rank features vs time different IDs

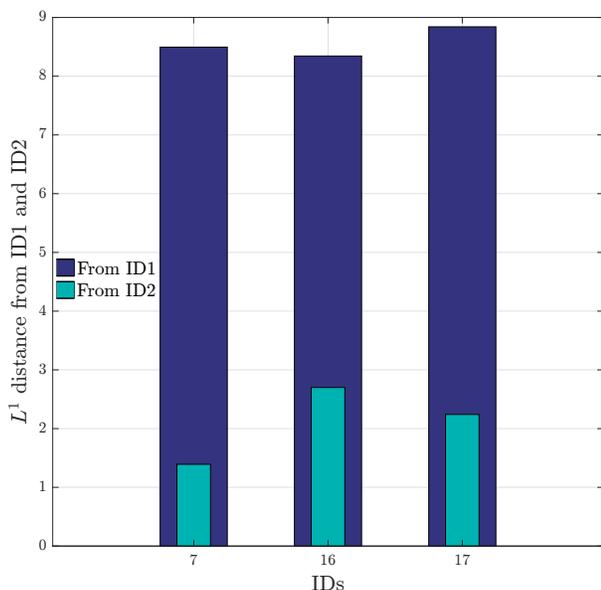


Fig. 21: Innovation vs time for various buffer sizes for the car video

C. Graph Edit Distance

To demonstrate the potential use of GED in a semantic extractor, we again take a simple classifier, and its confusion matrix given in Fig. 22. Note that such a simple classifier will only generate isolated nodes based on its detected classes. Therefore, the edit costs that need to be defined are the node insertion, deletion, and substitution costs, as listed in Table I. The substitution costs are defined as the probability of confusion between any two classes; and are shown in Fig. 23. Note that some entries show an infinite cost, i.e., meaning some specific patterns were never confused during the training. The insertion and deletion costs are calculated using (26) and Table I, and are shown in Fig. 24. These cost definitions can be used in (25) to compute GED metrics for any output this particular semantic extractor can generate.

To further showcase the application of the GED in tracking and filtering of semantic information, we introduce an

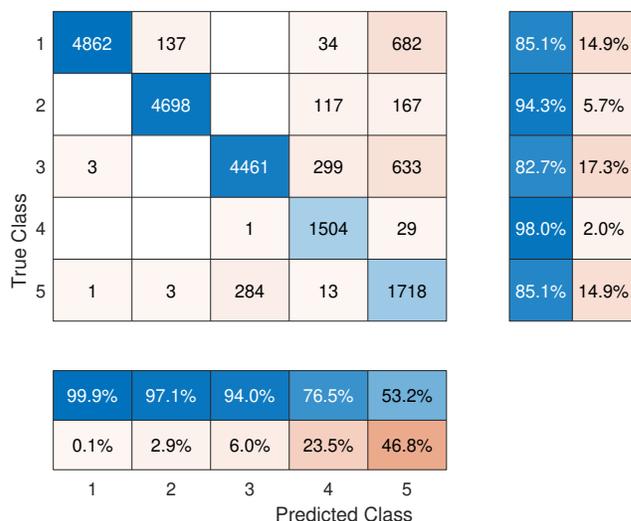


Fig. 22: Confusion matrix example for a classifier with five outputs with $\tau = 0.9$. The additional tables at the right and bottom side denote the accuracy rates of each row and column, respectively.

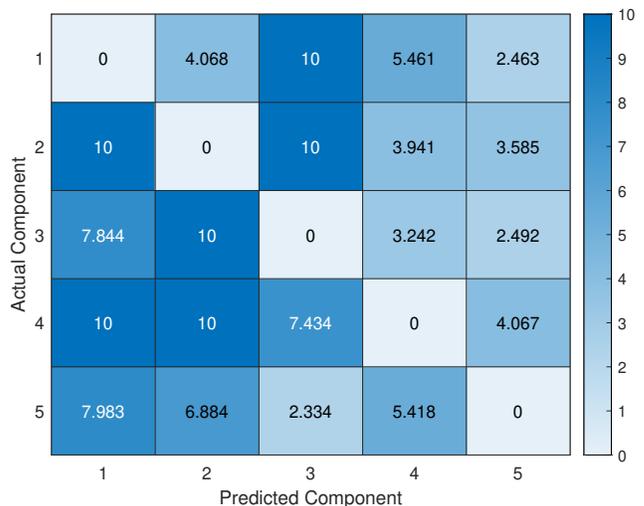


Fig. 23: Node substitution costs for the confusion matrix given in Fig. 22.

imperfect predicate detector that can identify three different types of relationships among the components detected by the classifier. The *conditional* confusion matrix and the corresponding GED costs for the predicate detector are calculated but not given here for brevity. Then, we simulate a temporal scenario given in Fig. 25, where some of these five classes and three predicates can be identified (sometimes erroneously). The generated scenario is then processed with the example classifier and predicate detector explained in this section. The resulting graph signals g_t at time t are compared to a baseline graph g_{base} , which is updated to be the latest detected graph with at least five consecutive detections. The expected and measured GED metrics with this configuration are presented in Fig. 26. As shown in Fig. 26a, each of the actual innovation events in the graph signal is identified with a jump in the

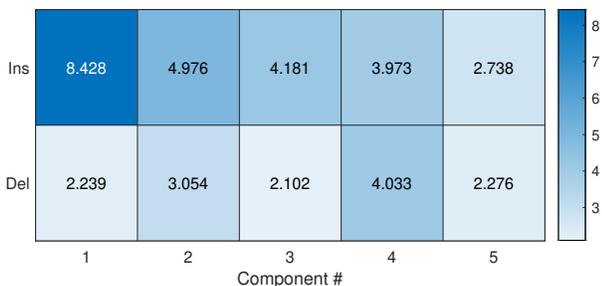


Fig. 24: Node insertion/deletion costs for the GED example.

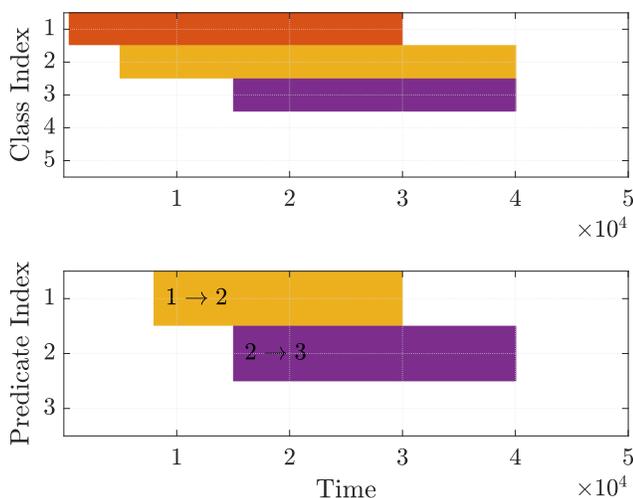


Fig. 25: Simulated scenario for the classifier and predicate outputs for 50,000 frames. The text annotations on the predicates indicate which two classes the predicate connects.

GED metric, followed by an updated baseline (return to zero) 5 frames after. Measured GEDs around different innovation time indices in Figs. 26b–26d show that sometimes there are erroneous detections that give rise to a nonzero GED; however, they can be easily filtered out by the straightforward algorithm explained above.

We also perform a GED-based pre-filtering on a video signal using the semantic extractor given in [9]. In this example, we have a 650-frame-long video (with 30fps) from a parking lot, where a car is maneuvering near several parked cars. This specific video is selected as its semantic output includes many false alarms, wrong categorizations, and missed detections. In particular, we are interested in showcasing the GED method’s ability to reconcile semantic confusions by exploiting the prior knowledge about the statistical similarity of the confused categories. A series of still images from the video example is shown in Fig. 27, with the semantic output summarized in Fig. 28.

As shown in Figs. 27 and 28, the maneuvering car (*Car-9*) is classified correctly as first, but then misclassified as two different boats as the video progresses. Focusing on this part of the video, as shown in Fig. 28b, we now demonstrate the use of GED in filtering out the semantic noise in this event. Firstly, we generate similar semantic outputs from a set of videos recorded by the same camera, in a similar time of the day to

TABLE II: Derived statistics for Boat and Car classes. Values on the left are the probabilities and the ones on the right are the corresponding GED costs.

	Car (observed)	Boat (observed)
Car (actual)	95% / 0.02	4.5% / 3.1
Prevalence	10% / 2.3	0.5% / 5.3

generate the confusion probabilities between the *Car* and *Boat* classes, as well as their prevalence (estimated rate of occurrence) in this particular setup. Using these statistics, the GED edit costs are calculated and shown in Table II. As shown in the top row of Table II, the confusion rate among the *Car* and the *Boat* classes is fairly small. However, taking into account the estimated prevalence of these classes and using Bayes’ rule, we can infer the posterior probability of a boat being confused by a car as $P(\text{Car exists}|\text{Boat observed}) = 90\%$ and the corresponding GED cost as $-\log(0.9) \approx 0.11$, which shows the statistical likelihood of a confusion is in fact very high. Using these cost definitions, the time evolution of the GED between consecutive semantic outputs is illustrated in Fig. 29. Note that by implementing a simple threshold, we can easily identify significant innovations and filter out the events that are not statistically significant. As shown in Fig. 29, by using a statistical significance threshold of 80%, or a corresponding GED cost limit of 0.2, we can reconcile the boat detections with the original correct identification. The GED output and the statistical significance limit are shown in Fig. 30.

As illustrated in Fig. 30, all misclassifications are correctly filtered except the duplicate detections over the same object. Since the GED method in its proposed form here only focuses on the time evolution of its scalar value, duplicate detections of the same object register as innovations. However, these are easily filtered with the next step in our proposed semantic extraction framework, with the attribute tracking module that can identify these detections to be from the same object.

D. Hidden Markov Model

In this part, HMM is applied to the raw semantic extraction of video signals to test its usefulness in practical scenarios. The main question is whether the smoothing algorithm (employing the Viterbi Algorithm) can correct the erroneously specified states when the model parameters ($\mathbf{A}, \mathbf{B}, \mathbf{p}$) are known. It is assumed that the scene starts with an empty graph; that is, there is not any component in the scene at the beginning. This is due to the design of the algorithm, which requires multiple consecutive detections of the same pattern to decide that it exists.

To illustrate the method, two video signals that contain cars and a person are used. For simplicity, a graph configuration with two predicates is employed. The predicates express the existence or absence of components in the graph language. So, the generated graph moves to a different state when a component enters or leaves the scene. The components leave or enter the scene independently of the other ones; hence if a scene contains N components, there are 2^N possible states in total. There are two classes of components, and the components can be connected to two predicates (can be in

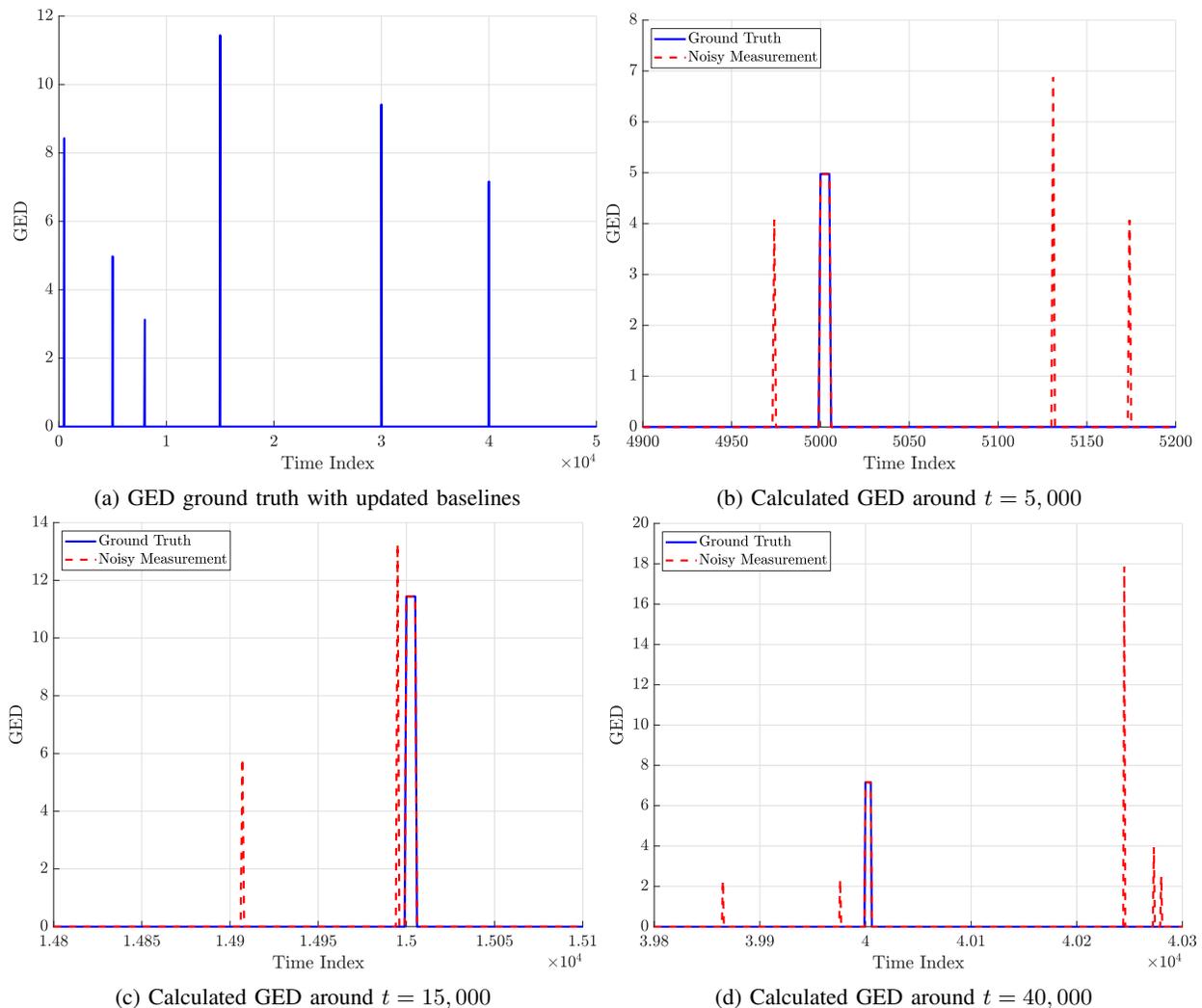


Fig. 26: GED example with updated baselines.

one of the two states). Assuming components of the same class follow the same statistical distribution, state transition matrices for car and person can be expressed as follows:

$$A_{car} = \begin{bmatrix} P_{0,0}^c & P_{0,1}^c(1 - P_{0,0}^c) \\ P_{1,0}^c(1 - P_{1,1}^c) & P_{1,1}^c \end{bmatrix},$$

$$A_{person} = \begin{bmatrix} P_{0,0}^p & P_{0,1}^p(1 - P_{0,0}^p) \\ P_{1,0}^p(1 - P_{1,1}^p) & P_{1,1}^p \end{bmatrix}.$$

The state transition matrix and observation kernel assumed for these videos are as follows:

$$A_{car} = \begin{bmatrix} 0.75 & 0.25 \\ 0.20 & 0.80 \end{bmatrix} \quad B_{car} = \begin{bmatrix} 0.75 & 0.25 \\ 0.40 & 0.60 \end{bmatrix}$$

$$A_{person} = \begin{bmatrix} 0.70 & 0.30 \\ 0.20 & 0.80 \end{bmatrix} \quad B_{person} = \begin{bmatrix} 0.70 & 0.30 \\ 0.40 & 0.60 \end{bmatrix}.$$

The true, observed, and estimated state sequences for the first video example are given in Fig. 31. The true state sequence is defined manually and represents the state sequence that would be observed if the semantic extractor did not

make any errors. The Viterbi Algorithm smooths out the errors introduced by the raw semantic extractor in the first video. As seen in Fig. 31, the semantic extractor erroneously decides that Car-2 is not in the scene in some instances. The Viterbi Algorithm corrects this mistake without introducing any further errors. For the second video illustrated in Fig. 32, the raw semantic extractor produces incorrect results for Car-1, Car-2, and Person-9 in some time instances. The proposed algorithm corrects these errors except for a short interval in Car-2's state sequence, where Car-2 is erroneously missing from the detections.

Note that in the second video, Car-2, Car-7, Car-16, and Car-17 are different instance identifications of the same underlying object (Car-2). The proposed HMM and Viterbi algorithm cannot reconcile these erroneous multiple identifications by itself. To produce a smooth and reliable output while reconciling among multiple identifications, we use the methods presented in this paper sequentially, as shown in Fig. 3. The *Attribute Tracking* example given in Section IX-B in Figs. 19–21 specifically illustrates the reconciliation among the multiple identifications of Car-2, Car-7, Car-16 and Car-17. Using the

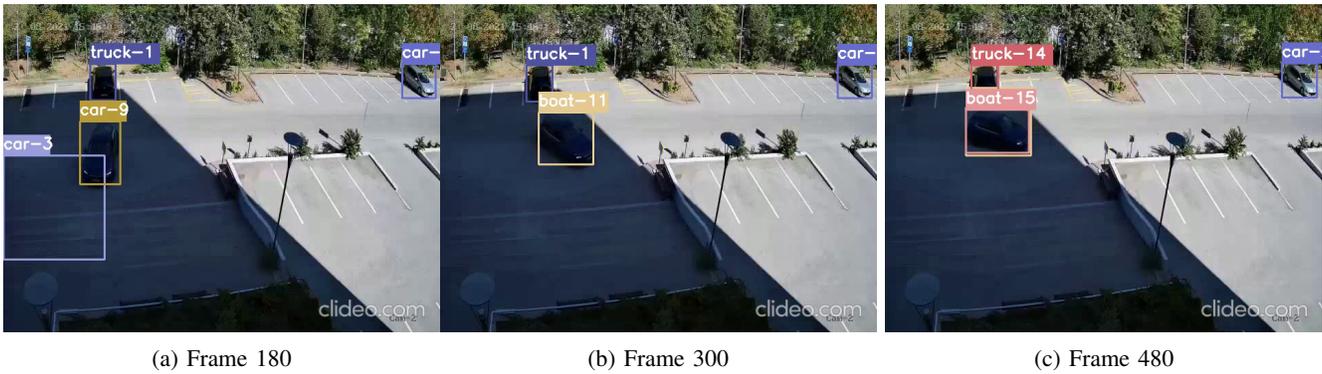


Fig. 27: Computer vision example with misdetections class outputs.

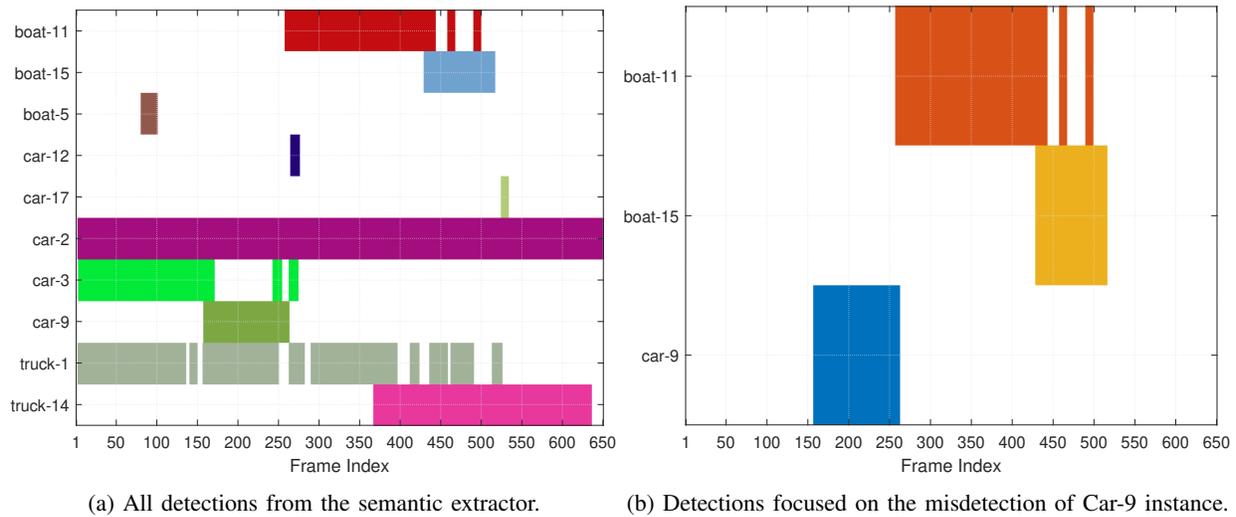


Fig. 28: Semantic extractor output with misdetections class outputs.

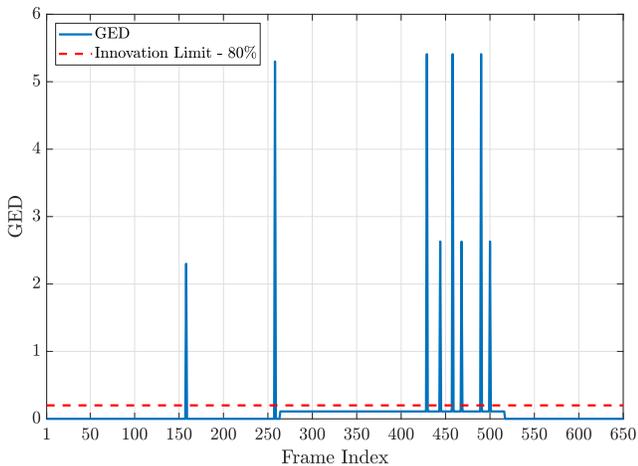


Fig. 29: GED across time for the events in Fig. 28.

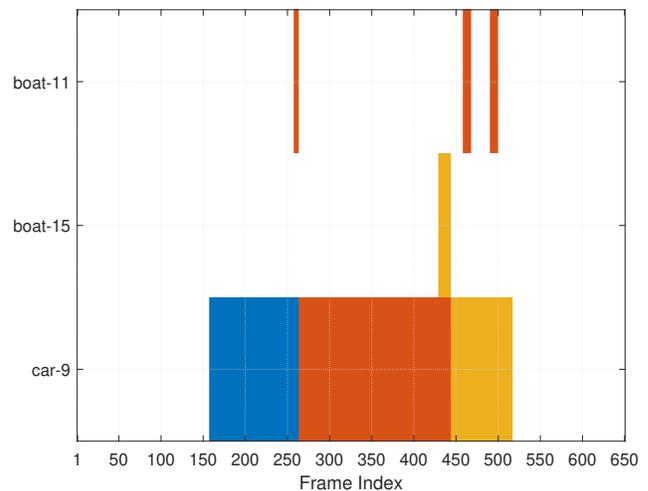


Fig. 30: Smoothed semantic evolution after GED filtering.

reconciled output of the *attribute tracking* module as the input of the *graph signal tracking* module, the output of the filtered semantic output is shown in Fig. 33. As illustrated in Fig. 33, the Viterbi algorithm is shown to correct the mistakes for each component, producing a filtered state sequence that is identical to the true sequence.

X. DISCUSSION ON THE SEMANTIC RATE OF INNOVATION

The proposed semantic extraction framework that has been defined and demonstrated so far is shown to generate reliable semantic graph signals with embedded numerical attributes. The proposed framework and its building blocks also enable

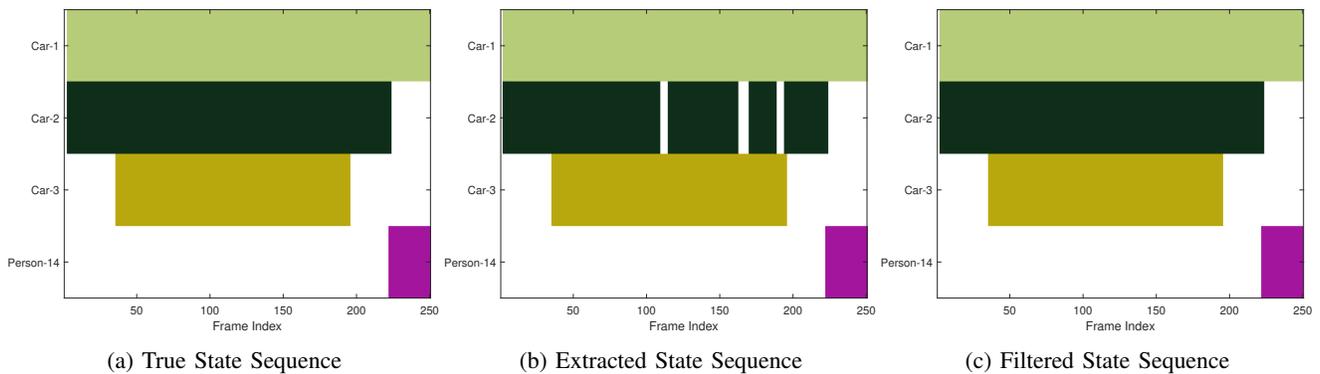


Fig. 31: State sequences for the first video. Semantic extractor labels Car-2 as absent when it is in the scene for some time instances. Viterbi Algorithm is able to correct the errors in this particular example.

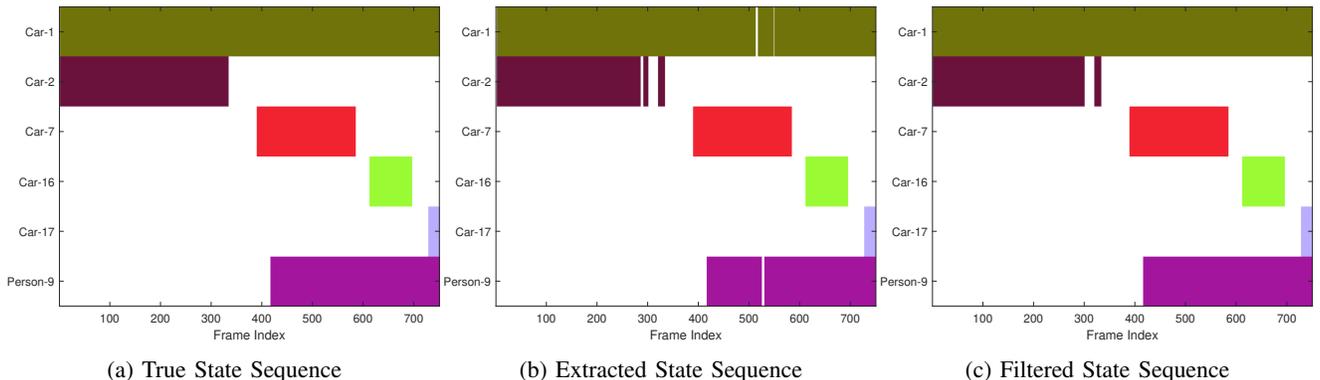


Fig. 32: State sequences for the second video. Semantic extractor labels Car-1, Car-2 and Person-9 as absent incorrectly in some time instances. Viterbi Algorithm is shown to corrects these mistakes except for a short blank interval for Car-2.

the identification of the innovation events either at the graph level (using HMM and/or GED) or at the attribute level (using subspace tracking). These events can be used to schedule transmission and/or storage events, depending on the device and the application of interest. Once the innovation events can be detected given a semantic signal, underlying statistical properties of innovations can be modeled and estimated. These estimates on innovation statistics can, in turn, be used to estimate the data throughput in a communications network.

Consider the semantic graph structure presented in Section III, with atomic bipartite graphs denoted by D_i that evolve across time as shown in Fig. 2. Let the rate of innovation at the graph level for D_i be I_{D_i} , and the corresponding average message length be d_{D_i} . The rates of innovation and the average message lengths at the attribute level can be similarly defined as

$$I_{A_i} = [I_{A_i}^1, I_{A_i}^2, \dots, I_{A_i}^L], \quad (27)$$

$$d_{A_i} = [d_{A_i}^1, d_{A_i}^2, \dots, d_{A_i}^L], \quad (28)$$

with L being the total number of attribute layers in (4). The total rate of transmitted or stored information for N atomic graphs without any goal-oriented filtering can be written as

$$R = \sum_{i=1}^N \left(I_{D_i} d_{D_i} + \sum_{l=1}^L I_{A_i} d_{A_i} \right). \quad (29)$$

Note that the semantic structure (components, predicates, and attributes) must be defined to optimize the rate of transmission using (29). Intuitively, the higher abstraction provided by the semantic structure should lead to lower rates of innovation for the graph and individual attributes compared to the rate of innovation of the raw input signal.

Building on (29), another significant advantage of the highly organized and hierarchical semantic structure is the processing of the graph signals in a goal-oriented manner. Given a dynamic goal \mathcal{G}_t (either internally or externally defined) that defines interest over a subset of the extracted semantic information, the rate of information can be significantly reduced. If we denote the most general goal that is interested in every output of the semantic extractor as \mathcal{G}_0 , with \mathcal{G}_t being a proper subset of \mathcal{G}_0 , application of the goal \mathcal{G}_t at a class and attribute level will lead to a reduced number of graphs $\hat{N} < N$ with a reduced number of attribute layers $\hat{L} < L$. The corresponding goal-oriented transmission rate can be written as

$$\hat{R} \triangleq \sum_{i=1}^{\hat{N}} \left(I_{D_i} d_{D_i} + \sum_{l=1}^{\hat{L}} I_{A_i} d_{A_i} \right), \quad (30)$$

with $\hat{R} < R$ for $\mathcal{G}_t \subset \mathcal{G}_0$.

Some of the most significant advantages of moving beyond technical communications toward goal-oriented semantic communications are the reduction of the rates of innovation

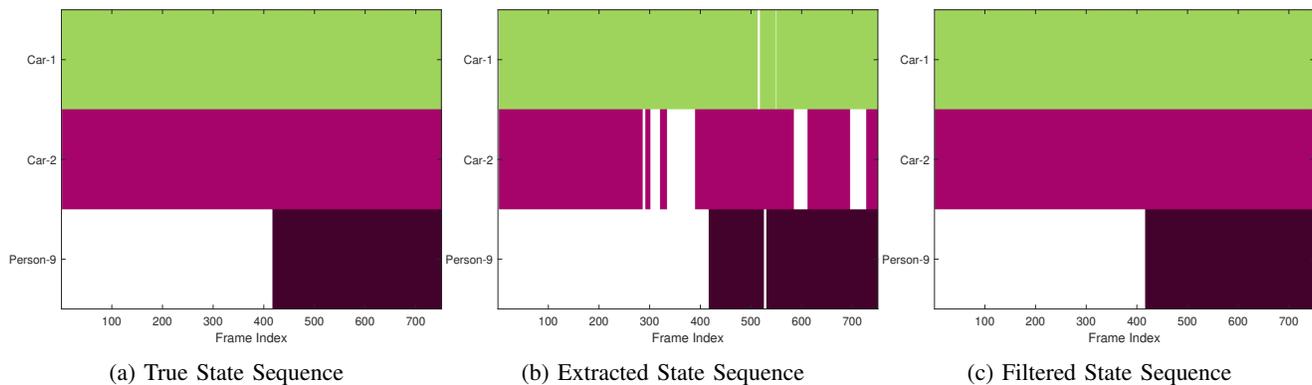


Fig. 33: State sequences for the second video with reconciliation of the components, combined state sequence of Car-2, 7, 16, 17 in 32 is represented as Car-2. Semantic extractor labels Car-1, Car-2 and Person-9 as absent incorrectly in some time instances, Viterbi algorithm is shown to correct these mistakes. The model parameters are selected as $P_{S_0, S_0} = 0.80$, $P_{S_1, S_1} = 0.85$, $P_{S_0, O_0} = 0.70$, $P_{S_1, O_1} = 0.55$ for the car class and $P_{S_0, S_0} = 0.90$, $P_{S_1, S_1} = 0.90$, $P_{S_0, O_0} = 0.65$, $P_{S_1, O_1} = 0.55$ for the person class.

and the inherent compression of the underlying signals. The modeling of the rates of innovation as well as the efficient compression of the semantic signals under dynamic goals will enable massive deployment for the next generation of sensor networks.

XI. CONCLUSIONS AND FUTURE RESEARCH DIRECTIONS

The next generation of signal processing and communication systems will employ intelligent agents that can generate semantic information from their local environment. This work introduces a semantic extraction framework, where the extracted graph-based imperfect semantic signals can be improved for better fidelity, filtered for semantic source noise, while enabling the identification of significant innovations in the semantic signal. The aforementioned tasks are achieved by exploiting known semantic characteristics of the environment and statistical characteristics of the front-end semantic extractor. The proposed methods provide reliable semantic outputs and enable efficient ways of identifying semantic innovation, while filtering out unwanted semantic noise. Note that the proposed metrics and methods can also be used to schedule transmission and storage events in semantics-enabled sensor devices.

As the semantic signal processing and the semantic extraction frameworks proposed in this work and in the literature move closer to massive deployment in the next generation of sensor networks, continued research on the practical implications of the proposed methods is required. Specifically, semantic extraction techniques for different signal modalities should be developed and standardized for a shared semantic structure/language for next-generation devices. Given standardized semantic extractors, semantically-aware time integration metrics should be modeled and estimated for application specific scenarios. Model estimation is also critical for the efficient implementation of the HMM-based modeling and processing proposed in this work.

The Markov model and its extension (HMM) also enable quantifying the semantic rate of innovation through the entropy

rate of the underlying Markov chain and the development of efficient compression and transmission schemes of semantic information. The modeling of the state sequences can further be improved by replacing the Markov model with a semi-Markov model, which can relieve the restriction of HMM on waiting time distributions and increase the representational power of the model.

We finally note that, depending on the type of sensor/device and its computational capabilities, the proposed methods can be used collectively or independently. As the signal processing and communications paradigms move towards semantic signal processing and transmission, we believe the proposed semantic extraction framework will be an essential building block in developing the next generation of sensor devices and networks.

REFERENCES

- [1] A. Rehman and T. Saba, "Features extraction for soccer video semantic analysis: current achievements and remaining issues," *Artificial Intelligence Review*, vol. 41, no. 3, pp. 451–461, 2014.
- [2] E. Çakır, G. Parascandolo, T. Heittola, H. Huttunen, and T. Virtanen, "Convolutional recurrent neural networks for polyphonic sound event detection," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 6, pp. 1291–1303, 2017.
- [3] A. Mesaros, T. Heittola, T. Virtanen, and M. D. Plumbley, "Sound event detection: A tutorial," *IEEE Signal Processing Magazine*, vol. 38, no. 5, pp. 67–83, 2021.
- [4] D. Cheng, F. Yang, S. Xiang, and J. Liu, "Financial time series forecasting with multi-modality graph neural network," *Pattern Recognition*, vol. 121, p. 108218, 2022.
- [5] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y.-J. A. Zhang, "The roadmap to 6G: AI empowered wireless networks," *IEEE Communications Magazine*, vol. 57, no. 8, pp. 84–90, 2019.
- [6] E. C. Strinati and S. Barbarossa, "6G networks: Beyond Shannon towards semantic and goal-oriented communications," *Computer Networks*, vol. 190, p. 107930, 2021.
- [7] I. F. Akyildiz, A. Kak, and S. Nie, "6G and beyond: The future of wireless communications systems," *IEEE Access*, vol. 8, pp. 133 995–134 030, 2020.
- [8] S. Dang, O. Amin, B. Shihada, and M.-S. Alouini, "What should 6G be?" *Nature Electronics*, vol. 3, no. 1, pp. 20–29, 2020.
- [9] M. Kalfa, M. Gok, A. Atalik, B. Tegin, T. M. Duman, and O. Arikan, "Towards goal-oriented semantic signal processing: Applications and future challenges," *Digital Signal Processing*, vol. 119, p. 103134, 2021.
- [10] G. Shi, D. Gao, X. Song, J. Chai, M. Yang, X. Xie, L. Li, and X. Li, "A new communication paradigm: from bit accuracy to semantic fidelity," *arXiv preprint arXiv:2101.12649*, 2021.

- [11] Q. Hu, G. Zhang, Z. Qin, Y. Cai, and G. Yu, "Robust semantic communications against semantic noise," *arXiv preprint arXiv:2202.03338*, 2022.
- [12] H. Xie, Z. Qin, G. Y. Li, and B.-H. Juang, "Deep learning enabled semantic communication systems," *IEEE Transactions on Signal Processing*, vol. 69, pp. 2663–2675, 2021.
- [13] C. E. Shannon, "A mathematical theory of communication," *The Bell system technical journal*, vol. 27, no. 3, pp. 379–423, 1948.
- [14] W. Weaver, "Recent contributions to the mathematical theory of communication," *ETC: a review of general semantics*, pp. 261–281, 1953.
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [16] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580–587.
- [17] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [18] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.
- [19] M. Tan, R. Pang, and Q. V. Le, "EfficientDet: Scalable and efficient object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10781–10790.
- [20] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Scaled-YOLOv4: Scaling cross stage partial network," in *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, 2021, pp. 13 029–13 038.
- [21] W. Van Gansbeke, S. Vandenhende, S. Georgoulis, and L. Van Gool, "Unsupervised semantic segmentation by contrasting object mask proposals," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10052–10062.
- [22] X. Xia and B. Kulis, "W-net: A deep model for fully unsupervised image segmentation," *arXiv preprint arXiv:1711.08506*, 2017.
- [23] X. Ji, J. F. Henriques, and A. Vedaldi, "Invariant information distillation for unsupervised image segmentation and clustering," *arXiv preprint arXiv:1807.06653*, vol. 2, no. 3, p. 8, 2018.
- [24] J. Shotton, M. Johnson, and R. Cipolla, "Semantic texton forests for image categorization and segmentation," in *2008 IEEE conference on computer vision and pattern recognition*. IEEE, 2008, pp. 1–8.
- [25] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," in *CVPR 2011*. IEEE, 2011, pp. 1297–1304.
- [26] J. Tighe, M. Niethammer, and S. Lazebnik, "Scene parsing with object instances and occlusion ordering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3748–3755.
- [27] S. Hao, Y. Zhou, and Y. Guo, "A brief survey on semantic segmentation with deep learning," *Neurocomputing*, vol. 406, pp. 302–321, 2020.
- [28] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [29] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [30] R. P. Poudel, S. Liwicki, and R. Cipolla, "Fast-SCNN: Fast semantic segmentation network," *arXiv preprint arXiv:1902.04502*, 2019.
- [31] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, "Panoptic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9404–9413.
- [32] F. Lateef and Y. Ruichek, "Survey on semantic segmentation using deep learning techniques," *Neurocomputing*, vol. 338, pp. 321–348, 2019.
- [33] A. Garcia-Garcia, S. Orts-Escobedo, S. Oprea, V. Villena-Martinez, P. Martinez-Gonzalez, and J. Garcia-Rodriguez, "A survey on deep learning techniques for image and video semantic segmentation," *Applied Soft Computing*, vol. 70, pp. 41–65, 2018.
- [34] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2625–2634.
- [35] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3128–3137.
- [36] J. Mao, W. Xu, Y. Yang, J. Wang, and A. L. Yuille, "Explain images with multimodal recurrent neural networks," *arXiv preprint arXiv:1410.1090*, 2014.
- [37] X. Chen and C. Lawrence Zitnick, "Mind's eye: A recurrent visual representation for image caption generation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2422–2431.
- [38] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International conference on machine learning*. PMLR, 2015, pp. 2048–2057.
- [39] J. Johnson, A. Karpathy, and L. Fei-Fei, "DenseCap: Fully convolutional localization networks for dense captioning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4565–4574.
- [40] L. Yang, K. Tang, J. Yang, and L.-J. Li, "Dense captioning with joint inference and visual context," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2193–2202.
- [41] D.-J. Kim, T.-H. Oh, J. Choi, and I. S. Kweon, "Dense relational image captioning via multi-task triple-stream networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [42] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International journal of computer vision*, vol. 123, no. 1, pp. 32–73, 2017.
- [43] J. Krause, J. Johnson, R. Krishna, and L. Fei-Fei, "A hierarchical approach for generating descriptive image paragraphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 317–325.
- [44] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, and J. Carlos Niebles, "Dense-captioning events in videos," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 706–715.
- [45] H. Xu, B. Li, V. Ramanishka, L. Sigal, and K. Saenko, "Joint event detection and description in continuous video streams," in *2019 IEEE winter conference on applications of computer vision (WACV)*. IEEE, 2019, pp. 396–405.
- [46] M. Z. Hossain, F. Sohel, M. F. Shiratuddin, and H. Laga, "A comprehensive survey of deep learning for image captioning," *ACM Computing Surveys (CSUR)*, vol. 51, no. 6, pp. 1–36, 2019.
- [47] N. Aafaq, A. Mian, W. Liu, S. Z. Gilani, and M. Shah, "Video description: A survey of methods, datasets, and evaluation metrics," *ACM Computing Surveys (CSUR)*, vol. 52, no. 6, pp. 1–37, 2019.
- [48] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, and L. Fei-Fei, "Image retrieval using scene graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3668–3678.
- [49] Y. Li, W. Ouyang, B. Zhou, K. Wang, and X. Wang, "Scene graph generation from objects, phrases and region captions," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1261–1270.
- [50] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei, "Scene graph generation by iterative message passing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5410–5419.
- [51] Y. Li, W. Ouyang, X. Wang, and X. Tang, "ViP-CNN: Visual phrase guided convolutional neural network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1347–1356.
- [52] Y. Li, W. Ouyang, B. Zhou, J. Shi, C. Zhang, and X. Wang, "Factorizable net: an efficient subgraph-based framework for scene graph generation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 335–351.
- [53] J. Yang, J. Lu, S. Lee, D. Batra, and D. Parikh, "Graph R-CNN for scene graph generation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 670–685.
- [54] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [55] K. Tang, H. Zhang, B. Wu, W. Luo, and W. Liu, "Learning to compose dynamic tree structures for visual contexts," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 6619–6628.
- [56] R. Zellers, M. Yatskar, S. Thomson, and Y. Choi, "Neural motifs: Scene graph parsing with global context," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5831–5840.
- [57] R. Wang, Z. Wei, P. Li, Q. Zhang, and X. Huang, "Storytelling from an image stream using scene graphs," in *Proceedings of the AAAI*

- Conference on Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 9185–9192.
- [58] H. Qi, Y. Xu, T. Yuan, T. Wu, and S.-C. Zhu, “Scene-centric joint parsing of cross-view videos,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [59] X. Chang, P. Ren, P. Xu, Z. Li, X. Chen, and A. Hauptmann, “Scene graphs: A survey of generations and applications,” *arXiv preprint arXiv:2104.01111*, 2021.
- [60] A. Agarwal, A. Mangal *et al.*, “Visual relationship detection using scene graphs: A survey,” *arXiv preprint arXiv:2005.08045*, 2020.
- [61] B. Güler, A. Yener, and A. Swami, “The semantic communication game,” *IEEE Transactions on Cognitive Communications and Networking*, vol. 4, no. 4, pp. 787–802, 2018.
- [62] Z. Qin, X. Tao, J. Lu, and G. Y. Li, “Semantic communications: Principles and challenges,” *arXiv preprint arXiv:2201.01389*, 2021.
- [63] G. Shi, Y. Xiao, Y. Li, and X. Xie, “From semantic communication to semantic-aware networking: Model, architecture, and open problems,” *IEEE Communications Magazine*, vol. 59, no. 8, pp. 44–50, 2021.
- [64] E. Uysal, O. Kaya, A. Ephremides, J. Gross, M. Codreanu, P. Popovski, M. Assaad, G. Liva, A. Munari, T. Soleymani *et al.*, “Semantic communications in networked systems,” *arXiv preprint arXiv:2103.05391*, 2021.
- [65] M. Kountouris and N. Pappas, “Semantics-empowered communication for networked intelligent systems,” *IEEE Communications Magazine*, vol. 59, no. 6, pp. 96–102, 2021.
- [66] Q. Lan, D. Wen, Z. Zhang, Q. Zeng, X. Chen, P. Popovski, and K. Huang, “What is semantic communication? a view on conveying meaning in the era of machine intelligence,” *Journal of Communications and Information Networks*, vol. 6, no. 4, pp. 336–371, 2021.
- [67] H. Xie and Z. Qin, “A lite distributed semantic communication system for internet of things,” *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 142–153, 2020.
- [68] Y. Wang, M. Chen, W. Saad, T. Luo, S. Cui, and H. V. Poor, “Performance optimization for semantic communications: An attention-based learning approach,” in *2021 IEEE Global Communications Conference (GLOBECOM)*, 2021, pp. 1–6.
- [69] Z. Weng and Z. Qin, “Semantic communication systems for speech transmission,” *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, pp. 2434–2444, 2021.
- [70] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [71] Y. Yang, C. Guo, F. Liu, C. Liu, L. Sun, Q. Sun, and J. Chen, “Semantic communications with artificial intelligence tasks: Reducing bandwidth requirements and improving artificial intelligence task performance,” *IEEE Industrial Electronics Magazine*, 2022.
- [72] H. Yoo, T. Jung, L. Dai, S. Kim, and C.-B. Chae, “Real-time semantic communications with a vision transformer,” *arXiv preprint arXiv:2205.03886*, 2022.
- [73] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [74] T.-Y. Tung, D. B. Kurka, M. Jankowski, and D. Gunduz, “DeepJSCC-Q: Constellation Constrained Deep Joint Source-Channel Coding,” *arXiv preprint arXiv:2206.08100*, 2022.
- [75] E. Boursoulatzé, D. B. Kurka, and D. Gündüz, “Deep joint source-channel coding for wireless image transmission,” *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 3, pp. 567–579, 2019.
- [76] T.-Y. Tung and D. Gündüz, “DeepWiVe: Deep-learning-aided wireless video transmission,” *arXiv preprint arXiv:2111.13034*, 2021.
- [77] T.-Y. Tung, S. Kobus, J. P. Roig, and D. Gündüz, “Effective communications: A joint learning and communication framework for multi-agent reinforcement learning over noisy channels,” *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, pp. 2590–2603, 2021.
- [78] J. Bao, P. Basu, M. Dean, C. Partridge, A. Swami, W. Leland, and J. A. Hendler, “Towards a theory of semantic communication,” in *Proceedings of the 2011 IEEE 1st International Network Science Workshop*, 2011, pp. 110–117.
- [79] E. Che, J. Jung, and M. J. Olsen, “Object recognition, segmentation, and classification of mobile laser scanning point clouds: A state of the art review,” *Sensors*, vol. 19, no. 4, 2019.
- [80] G. Dillard and J. Rickard, “Performance of an MTI followed by incoherent integration for nonfluctuating signals,” in *International Radar Conference*, 1980, pp. 194–199.
- [81] X. Li, G. Cui, W. Yi, and L. Kong, “Coherent integration for maneuvering target detection based on radon-lv’s distribution,” *IEEE Signal Processing Letters*, vol. 22, no. 9, pp. 1467–1471, 2015.
- [82] X. Li, L. Kong, G. Cui, and W. Yi, “CLEAN-based coherent integration method for high-speed multi-targets detection,” *IET Radar, Sonar & Navigation*, vol. 10, no. 9, pp. 1671–1682, 2016.
- [83] S. Visa, B. Ramsay, A. L. Ralescu, and E. Van Der Knaap, “Confusion matrix-based feature selection,” *MAICS*, vol. 710, pp. 120–127, 2011.
- [84] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, “3D object representations for fine-grained categorization,” in *4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13)*, Sydney, Australia, 2013.
- [85] N. Wojke, A. Bewley, and D. Paulus, “Simple online and realtime tracking with a deep association metric,” in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 3645–3649.
- [86] N. Aharon, R. Orfaig, and B.-Z. Bobrovsky, “BoT-SORT: Robust associations multi-pedestrian tracking,” *arXiv preprint arXiv:2206.14651*, 2022.
- [87] Y. Zhang, P. Sun, Y. Jiang, D. Yu, Z. Yuan, P. Luo, W. Liu, and X. Wang, “ByteTrack: Multi-object tracking by associating every detection box,” *arXiv preprint arXiv:2110.06864*, 2021.
- [88] Y. Du, Y. Song, B. Yang, and Y. Zhao, “StrongSORT: Make DeepSORT great again,” *arXiv preprint arXiv:2202.13514*, 2022.
- [89] N. Vaswani, T. Bouwmans, S. Javed, and P. Narayanamurthy, “Robust subspace learning: Robust PCA, robust subspace tracking, and robust subspace recovery,” *IEEE Signal Processing Magazine*, vol. 35, no. 4, pp. 32–55, 2018.
- [90] P. Rodríguez and B. Wohlberg, “Fast principal component pursuit via alternating minimization,” in *2013 IEEE International Conference on Image Processing*, 2013, pp. 69–73.
- [91] L. Rabiner, “A tutorial on hidden Markov models and selected applications in speech recognition,” *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [92] A. Viterbi, “Error bounds for convolutional codes and an asymptotically optimum decoding algorithm,” *IEEE Transactions on Information Theory*, vol. 13, no. 2, pp. 260–269, 1967.
- [93] L. E. Baum and J. A. Eagon, “An inequality with applications to statistical estimation for probabilistic functions of Markov processes and to a model for ecology,” *Bulletin of the American Mathematical Society*, vol. 73, no. 3, pp. 360–363, 1967.
- [94] J. Anderson, “Limited search trellis decoding of convolutional codes,” *IEEE Transactions on Information Theory*, vol. 35, no. 5, pp. 944–955, 1989.
- [95] L. E. Baum and T. Petrie, “Statistical inference for probabilistic functions of finite state Markov chains,” *The Annals of Mathematical Statistics*, vol. 37, no. 6, pp. 1554–1563, 1966.
- [96] S.-Z. Yu, “Hidden semi-Markov models,” *Artificial Intelligence*, vol. 174, no. 2, pp. 215–243, 2010, special Review Issue.
- [97] A. Sanfeliu and K.-S. Fu, “A distance measure between attributed relational graphs for pattern recognition,” *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-13, no. 3, pp. 353–362, 1983.
- [98] X. Gao, B. Xiao, D. Tao, and X. Li, “A survey of graph edit distance,” *Pattern Analysis and applications*, vol. 13, no. 1, pp. 113–129, 2010.
- [99] D. B. Blumenthal and J. Gamper, “On the exact computation of the graph edit distance,” *Pattern Recognition Letters*, vol. 134, pp. 46–57, 2020.