# Towards Goal-Oriented Semantic Signal Processing: Applications and Future Challenges

Mert Kalfa[a,*], Mehmetcan Gok[a], Arda Atalik[a], Busra Tegin[a], Tolga M. Duman[a], Orhan Arikan[a]

*aDepartment of Electrical and Electronics Engineering, Bilkent University, 06800, Ankara, Turkey*

## Abstract

Advances in machine learning technology have enabled real-time extraction of semantic information in signals which can revolutionize signal processing techniques and improve their performance significantly for the next generation of applications. With the objective of a concrete representation and efficient processing of the semantic information, we propose and demonstrate a formal graph-based semantic language and a goal filtering method that enables goal-oriented signal processing. The proposed semantic signal processing framework can easily be tailored for specific applications and goals in a diverse range of signal processing applications. To illustrate its wide range of applicability, we investigate several use cases and provide details on how the proposed goal-oriented semantic signal processing framework can be customized. We also investigate and propose techniques for communications where sensor data is semantically processed and semantic information is exchanged across a sensor network.

*Keywords:* Semantic signal processing, semantic languages, semantic filtering, goal-oriented filtering, graph-based languages, semantic communications, goal-oriented communications.

## 1. Introduction

Signal processing has played a major role in the era of the digital revolution. It has been essential in the design and development of complex systems that integrate the contributions of multiple scientific disciplines. Significant challenges in critical areas including remote sensing, healthcare, finance, transportation, entertainment, communications, and security have been successfully overcome with advances in signal processing techniques [1].

Semantic information theory (SIT) [2] has been around almost as long as the classical information theory (CIT) [3]. As its name implies, SIT is focused on the *semantic problem* which is the accurate generation, processing, storage, and transmission of the *intended meaning* rather than individual bits and symbols, which is referred to as the *technical problem*. Despite its long history, semantic information has been used sporadically and only at a basic level in various signal processing applications so far. An important example of semantic information utilization is the target tracking by an individual sensor or by a sensor network. For instance, in computer vision applications, following the detection of objects via real-time processing of images or videos, typically, only a small subset of detected object reports are transferred to the tracking system that initiates new tracks or maintains the existing ones. In such processes, objects that are moving closely can be identified and tracked as a group. In this application, detection and tracking of objects with their identified group properties can be viewed as a semantic extraction of information from image or video signals. Furthermore, if multiple sensors operate as part of a network, they can share the semantic information in the form of object tracks for improved surveillance and object handover operations [4–6].

---

*Corresponding author

*Email addresses:* kalfa@ee.bilkent.edu.tr (Mert Kalfa), mehmetcan@ee.bilkent.edu.tr (Mehmetcan Gok), arda.atalik@bilkent.edu.tr (Arda Atalik), btegin@ee.bilkent.edu.tr (Busra Tegin), duman@ee.bilkent.edu.tr (Tolga M. Duman), oarikan@ee.bilkent.edu.tr (Orhan Arikan)

Recent advances in machine learning have far-reaching implications in a diverse range of disciplines including but not limited to healthcare [7–9], finance [10–12], entertainment [13, 14], and communications [15, 16]. Among many achievements of machine learning technology, real-time extraction of rich semantic information in streaming data is of critical importance in future signal processing applications. For instance, the identification and classification of objects in images and video streams can now be performed in real time [17–20]. Availability of such rich semantic data can significantly change the signal processing techniques and improve their performance.

The 6G and beyond communication systems have the potential to bring about an *Internet of Intelligence* with connected people, connected things, and connected intelligence [21–25]. It will possess the potential to build up a wireless network of *everything sensing*, *everything connected* and *everything intelligent*. Besides the traditional radio technology innovations, 6G and beyond will explore many innovative architectures and technologies, among which the goal-oriented semantic signal processing and communications will be a core component. Semantic communications will enable effective utilization of the available network capacity. Furthermore, the users of the network such as autonomous driving cars will benefit significantly from the goal filtered semantic data exchanges among themselves and the smart city backbone.

In this paper, for a structured and universal representation and efficient processing of the semantic information, we propose a formal semantic signal processing framework that includes a graph-based semantic language and a goal-oriented parsing method. The proposed framework can easily be tailored for specific signal processing applications and goals. In the proposed framework, a semantic information extractor identifies and classifies components of a signal into a set of application-specific classes. The semantic relations among the identified components are described by a set of application-specific predicates. Furthermore, along with the identification and classification of the components, each node in the graph is associated with a hierarchical set of attributes that provides additional information, ordered in increasing detail. Those components that are semantically related to each other form directed bipartite graphs where only edges between a component and a predicate are allowed to exist. Since typically there are very few nodes in these graphs, operations on these graphs can be implemented very efficiently. Furthermore, the proposed bipartite structure enables a complete yet relatively simple representation of signals and reduces the computational complexity of graph-based signal processing applications considerably.

In the proposed semantic signal processing framework, based on an internally or externally defined set of goals that can vary with time, graphs can be grouped as those which will be processed further and those that are currently not of interest. The further processing stages may include spatio-temporal tracking of graph parameters and conducting various operations on their attributes. At any point in the processing chain, the desired level of semantic information of those graphs which are of interest can be locally stored or shared with another processor through appropriate communication protocols. Since typical high bandwidth sensor data have a sparse occurrence of interesting events, the corresponding semantic signal processing algorithms result in remarkable compression rates. To illustrate the significant reduction in the communication rate compared to classical approaches, we investigate and propose techniques for communication over a sensor network where sensor data is semantically processed and semantic information is exchanged in the network. Furthermore, to demonstrate the wide range of applicability of the proposed goal-oriented semantic signal processing framework, we also investigate several use cases and provide details on how the proposed framework can be adapted for different signal modalities.

The rest of the paper is organized as follows. In the next section, we review existing approaches to semantic and goal-oriented information in signals. In Section 3, we present a detailed literature survey on the recently developed semantic information extraction and semantic communication techniques. The proposed goal-oriented semantic language and signal processing framework is introduced in Section 4. In Section 5, we demonstrate how the proposed semantic signal processing framework can be customized for a few use cases. In Section 6, we investigate the transmission and storage of semantic information, as well as the data compression potential of the proposed framework. Concluding remarks and promising research directions are given in Section 7.

## 2. Semantic and Goal-Oriented Information

In his seminal paper [3], Claude E. Shannon introduced CIT and paved the way for modern communications as we know it. However, as recognized and later formalized by Shannon and Weaver [26], CIT only addresses the low-level *technical problem*, i.e., only the accurate representation and transmission of bits and symbols are taken into account. At higher levels of abstraction, there are the *semantic problem* and the *effectiveness problem*, which focus on transmission of the *intended meaning* and execution of the *intended effect*, respectively. This work focuses on the mid-level *semantic problem*.

Consider the speech communication between two people. We use words as part of a shared language to convey certain information to one another. If person-A describes a visual scene to person-B, the exact image that forms in the mind of person-B is almost certainly not the same as that in person-A's mind. However, if the *intended* information is correctly conveyed, no semantic information is lost during transmission. Moreover, from a throughput perspective, a complex scene with a few meaningful components can be conveyed with a comparable number of sentences, greatly reducing the number of *symbols* that are transmitted compared to a pixel-by-pixel breakdown of the scene. Another important factor in human communications is the ability to ask questions, which restricts the domain of knowledge considerably when conveying information. In this work, we refer to questions that can be asked internally or externally as *goals*. Hence, the term *goal-oriented information* is used throughout the paper as semantic information that is filtered by pertinent questions.

With the introduction of the Internet of Things (IoT), the next-generation communication networks will be dominated by massive machine-type communications (mMTC) [22]. Introducing a *language* as a medium to describe signals and form queries about them enables a mode of communication between machines that is similar to humans and can greatly reduce the overall throughput and increase efficiency. Although this increase in efficiency is readily expected and has been known almost since the introduction of CIT [26, 2], efficient extraction of semantic information was not possible until the recent advances in machine learning (ML) and artificial intelligence (AI) technologies.

ML and AI techniques enable the development of novel signal processing algorithms and systems that can extract and use the semantic information in their inputs in real-time. Modern ML methods including convolutional and recurrent deep neural network (DNN) architectures [27–29], Long-Short Term Memory (LSTM) networks [30], and more recently, scene graph generation techniques [31–36] make it possible to efficiently extract semantic information in signals of widely different modality such as speech, image and video signals. The resulting semantic signals allow for the next generation of signal processing techniques to be implemented at a semantic level with a goal-oriented approach in real-time.

Unlike the classical signal processing approaches, in semantic signal processing, goals are expressed in a semantic language. To be able to define suitable performance metrics, and implement compression/coding schemes based on conveyed or inferred meaning of transmission, one needs to define a *language* that maps these meanings to a predefined syntactic structure. Therefore, it is critical to establish a semantic information and language model that is sufficiently general to be suitable in various signal processing applications, yet simple enough to be used by low-end agents with stringent power and computing limitations. In the rest of this section, a brief review of previously proposed semantic language modalities is presented.

### 2.1. Natural Languages

The most popular and intuitive candidates for a shared language in signal processing and mMTC applications are Natural Languages (NL). There is a vast amount of work on Natural Language Processing (NLP) [37] for the generation of NL sentences given an input signal, especially for human-machine interface applications. Examples include question answering [38–40], captioning of images and videos [41–44], and discourse parsing [45–47]. Relatively recently, there has been an increasing focus on semantic communications using NL as a basis as well [48–51].

NLs are attractive since they offer a universal language for all agents, regardless of their specific capabilities. However, this universality also comes with the requirement of processing a massive knowledge base (e.g., the English language) with inherent ambiguities and inconsistencies. This makes the NL approaches

unnecessarily complex for simple IoT sensors and similar machine-type applications. For such cases, an unambiguous and simple language structure is required.

### 2.2. Propositional Logic

The seminal work on semantic communications by Carnap and Bar-Hillel [2] introduces propositional logic as a semantic language. In a propositional logic language, the information is encoded in Boolean symbols corresponding to *components* and *primitive properties*. Common logical operations such as NOT ($\neg$), AND ($\wedge$), OR ($\vee$), IMPLIES ($\Rightarrow$) etc., are used to combine several symbols to form *molecular sentences* describing a state (e.g., $\neg(A \Rightarrow B) \wedge C$). Goals in this language can be formed similarly using propositional logic, effectively parsing the existing truth table of symbols for the desired pattern.

The propositional logic as a semantic language is attractive, since it can be tailored for specific applications and as a result, may not suffer from the complexity and ambiguity of NLs. However, it is challenging to incorporate numerical attributes such as position and velocity into the language; therefore, the resulting descriptions may lack complete information about the signals of interest.

### 2.3. Graph-Based Languages

Graph-based languages offer significant advantages over NLs as they are mathematical constructs that can represent components in a signal, as well as their relationships and states. Most popular applications of graph-based languages include scene graph generations from images and videos [31–36], knowledge graphs and graph-based question answering [52–55], and semantic web applications [56–58]. The question answering problems use sub-graphs called *motifs* to search for a pattern within a graph representation, enabling a goal-oriented approach [59–62]. Moreover, graph nodes and edges can include additional attributes to convey a more complete description of the underlying scene [63].

With the aforementioned attractive qualities, we believe that graphs with attributes can provide a structured and complete description for a variety of signals of interest. However, the majority of graph-based language models still depend on NLs as their backbone; hence, they can still be unnecessarily complex for simple machine-type applications with stringent efficiency requirements. Therefore, we advocate the use of a graph-based language that can be tailored to specific signal domains and applications of interest with a unique and relatively small knowledge base that would be of great advantage for machine-type applications.

## 3. Survey of Semantic Transformations and Goal-Oriented Semantic Communications

Before rigorously defining our proposed language and signal processing framework in Section 4, we present the state-of-the-art in existing semantic transformations that can effectively process signals of different modalities and goal-oriented semantic communications in this section.

### 3.1. State-of-the-Art in Semantic Transformations

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 as discussed in Section 2.

In this part, we give a detailed survey of the state-of-the-art semantic transformations. It is important to point out that although there are many input-target modalities, we only discuss a small, albeit an intuitive subset of transformation modalities mostly focused on visual input types.
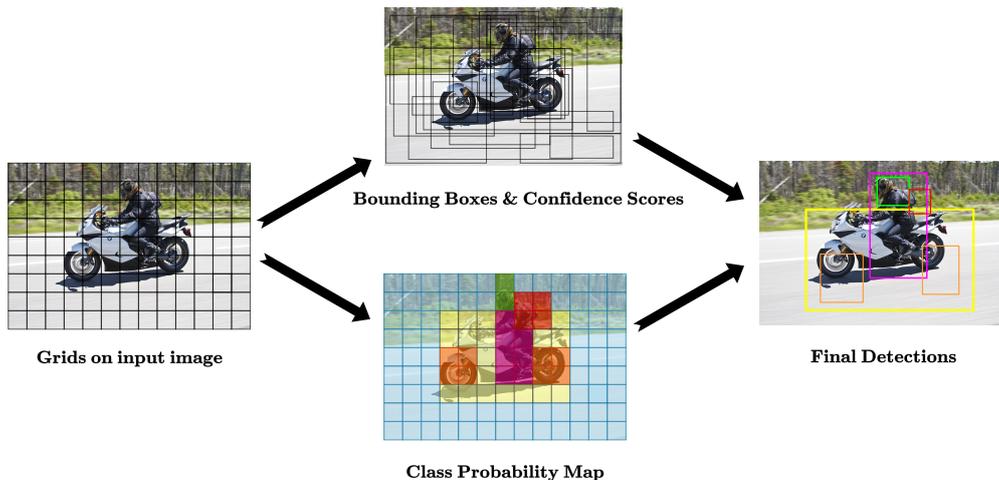
4

Figure 1: YOLO [17] object detection model. Unlike Faster-RCNN [67], YOLO does not use region proposals, instead, it splits the input image into grids and performs inference on each grid.

### 3.1.1. Object Detection and Segmentation

Object detection can be interpreted as one of the core semantic transformations from visual domain to domain of object classes. Convolutional Neural Networks (CNNs) [64] constitute fundamental processing modules of object detection methods. Recurrent CNN (RCNN) models are introduced in [65], where the selective search is utilized to propose candidate regions and each region is processed independently by a CNN and classified using support vector machines (SVMs). The use of RCNN models is further diversified in [66, 67]. In [66] instead of extracting region features separately, a single forward pass through a CNN is performed and the resulting feature map is branched into regions using Region-of-Interest (RoI) pooling. RCNNs with real-time processing capabilities are introduced in [67], where object regions are proposed by Region Proposal Networks (RPNs) instead of selective search. YOLO is introduced in [17] as the very first attempt on fast object detection task and it is further improved in [68, 69]. Unlike Faster-RCNN [67], YOLO, and its later versions, do not utilize region proposals. Instead, they split the input image into cells and perform inference on a limited number of boxes in each cell [17] as shown in Fig. 1. Moreover, recent works [70, 18] use network scaling approaches to achieve higher detection accuracy within shorter inference times, and obtain state-of-the-art results for real-time object detection. Particularly, [70] proposes a weighted bi-directional feature pyramid network to fuse features from multiple levels, and in [18] a network scaling method is applied to the YOLO architecture.

Segmentation is another semantic transformation that is applicable for visual inputs or inputs that are transformed to 2-D domains such as time-frequency or time-scale. Semantic segmentation can be considered as a higher resolution version of object detection where a category label is assigned to each pixel. More formally, semantic segmentation aims to find a way to assign a label from categories to a set of random variables that correspond to each pixel in the image. Here, each label can represent an object class such as "person", "plane", "car", etc., or labels can be the distinct but unspecified clusters in an unsupervised setting [71–73]. In [74], texton forests are proposed as efficient low-level features for image segmentation. Alternative approaches include random forest classifiers [75] and combination of SVMs and Markov Random Fields (MRFs) [76]. The state-of-the-art in semantic segmentation typically employs convolutional architectures in supervised, semi-supervised, and weakly-supervised settings [77]. It is important to note that, the features extracted by the deeper layers of a CNN are more concentrated on concise semantics with low spatial details whereas shallow layers of a CNN are more aware of spatial details such as edge orientations. Convolutional architectures for image segmentation include dilated convolutions [78] to incorporate context information from multiple scales, Fully Convolutional Network (FCN) architectures [79] using skip-connections [80], a symmetrical encoder-decoder structure called DeconvNet [81], and its simpler version in [82].
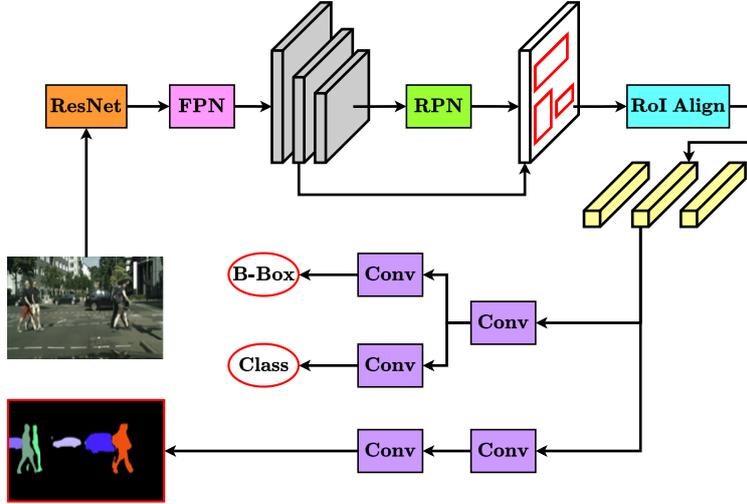
5

Figure 2: Mask-RCNN [90] model for instance segmentation. As its backbone, Mask-RCNN uses ResNet which is followed by a Feature Pyramid Network (FPN) and a Region Proposal Network (RPN). Features for proposed regions are extracted with RoI Alignment. Finally, bounding-box regression, instance classification, and segmentation mask inference are performed.

Recurrent architectures are also employed in semantic segmentation. In [83], Recurrent Fully Convolutional Networks are proposed for multi-slice Magnetic Resonance Imaging (MRI) segmentation where Gated Recurrent Unit (GRU) is incorporated into the bottleneck of the U-Net architecture given in [80]. Moreover, several Generative Adversarial Network (GAN) architectures are proposed [84–86], where adversarial training is introduced to semantic segmentation in [84].

Semantic segmentation architectures that offer real-time operation are proposed in [87–89]. In [90], Faster-RCNN [67] is modified for *instance segmentation* and Mask-RCNN is introduced, the process of which is illustrated in Fig. 2. Unlike semantic segmentation, instance segmentation aims at assigning labels to pixels at an object level as opposed to a class level. Union of instance and semantic level segmentation is called panoptic segmentation [91], where each pixel is associated with both instance and class level labels, as shown in Fig. 3 along with different segmentation types. Readers interested in semantic segmentation can find detailed taxonomy of models, applications, discussion on challenges, and possible future directions in [77, 92, 93].

### 3.1.2. Image and Video Captioning

An intuitive way of representing the semantic information embedded in images or videos is to describe them via natural languages (NLs). Image caption or annotation generation is defined as the process of generating textual descriptions for images, which include not only the descriptions of objects within frames but also their interrelations and states. Typical fundamental building blocks of captioning models include CNNs and Recurrent Neural Networks (RNNs). More specifically, a CNN backbone is utilized to extract the visual features and RNNs are used for sequence modeling [94–96, 94, 97–99]. To improve the caption quality, visual attention on CNNs [98] or captioning on multiple image regions are proposed [95]. In [99], dense captioning is introduced where object detection and caption generation tasks are tackled jointly in such a way that the detected visual concepts are described with short NL phrases. To overcome localization issues of visual concepts due to overlapping target regions, global image features are fused with region features for dense captioning in [100]. In [101], the dense captioning task is extended to relational captioning where multiple captions are generated for each object pair based on spatial, attentive, and contact relation information [102]. Examples of image captioning techniques are shown in Fig. 4.

Video captioning can be considered as a temporal extension of image captioning. Just like image captioning, the video captioning architectures also use CNNs (2D or 3D) as attention mechanisms or to extract visual semantic content. Then, RNNs or LSTMs are typically used to generate NL text sequences [103]

6

(a) Object Detection



(b) Semantic Segmentation



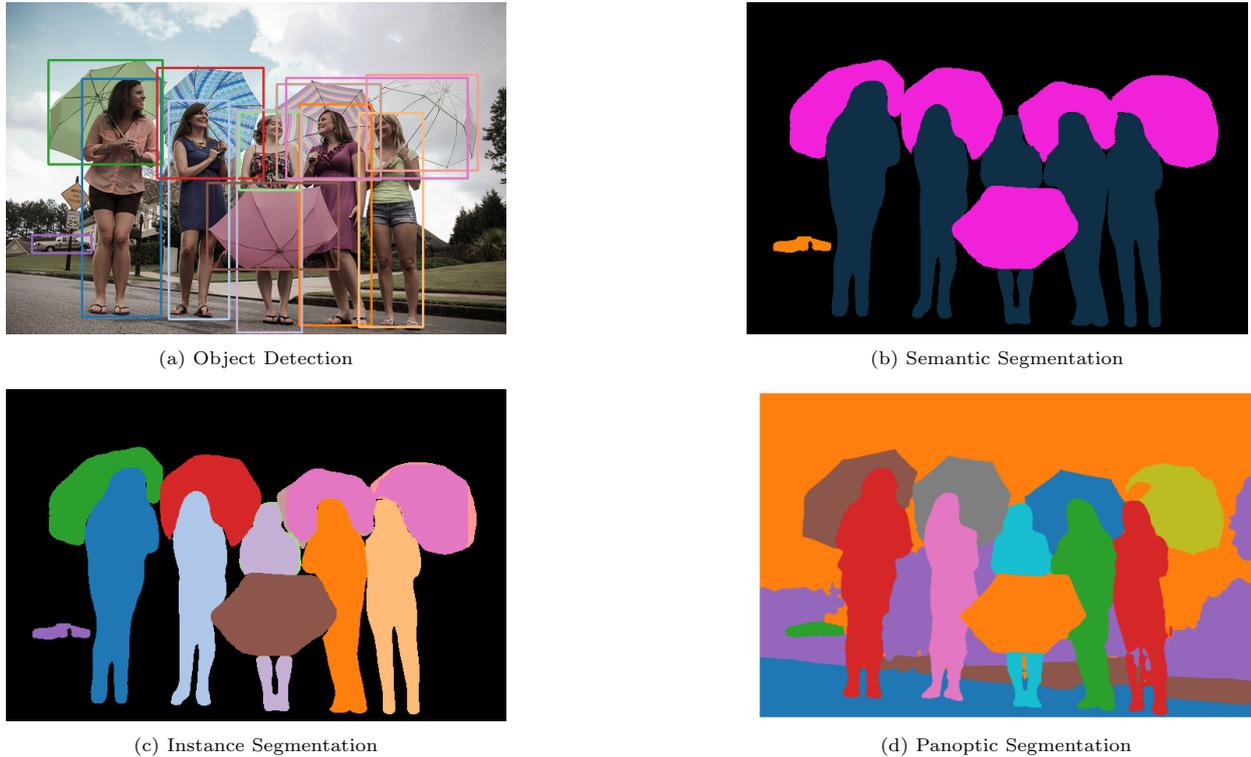(c) Instance Segmentation



(d) Panoptic Segmentation

Figure 3: Different object detection and segmentation techniques applied to an image. (a) In object detection, detected/classified objects are localized with bounding boxes. (b) Semantic segmentation aims to assign class labels at the pixel level. (c) Instance segmentation is applied to the object detection for finer localization. (d) In panoptic segmentation, both background and foreground are segmented such that each pixel is associated with an instance and a class simultaneously.

as illustrated in Fig. 5. Specifically, one of the early works [104], which is only applicable for videos of short duration, employs mean-pooling to frame representations extracted by a shared CNN and utilizes an LSTM architecture for caption generation. To extend the validity of extracted features to longer durations, recurrent visual encoder architectures are used [94, 105, 106]. In [107], instead of a single caption, multiple captions are generated without their temporal localizations using a hierarchical-RNN architecture. Furthermore, a hierarchical-RNN architecture is used to generate paragraphs for videos in [108]. Dense captioning works in the image domain are also extended to video signals in [44] where multiple captions are generated for each detected event, and they are temporally localized. Dense video captioning is also referred to as joint event detection and description generation. In [44] CD3 features are extracted beforehand, and they are fed into a proposal module that uses an attention mechanism. JEDDi-Net proposed in [109] is employed for dense video captioning in an end-to-end fashion without pre-feature extraction. It uses a 3D-CNN to extract the video features and a segment-proposal-network (SPN) to generate candidate segments for events. A comprehensive survey of image and video captioning models can be found in [110] and [103], respectively.

### 3.1.3. Scene Graph Generation

A powerful form of semantic transformation is to convert images into graphs that represent a scene and encode the visual relationships presented in the image. Scene graphs are proposed in [111] describe image features and object relationships in an explicit and structured way for image retrieval. Scene graph generation models can capture a higher-level of understanding of the scenes compared to object detection models by additionally identifying object attributes and relationships among them. In essence, scene graph generator models can be considered as image captioning models where they produce parsed versions of
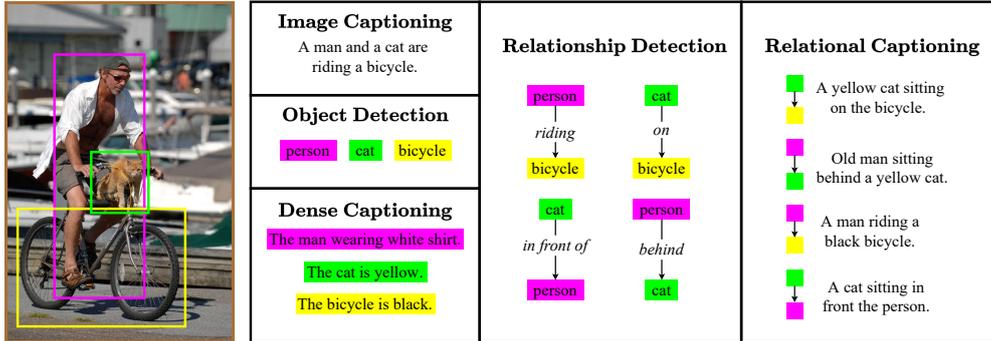
Figure 4: Some semantic transformation examples on the image domain. Image captioning simply generates a textual description of the image. Object detection identifies the objects residing in the image. Dense captioning generates a caption for each salient region or detected object in the image. Relationship detection identifies relations among each object pair. Relational captioning generates a caption for each detected object pair.

succinct captions that are represented in the graph format, instead of NL sentences. Specifically, a scene graph is a graphical data structure that describes the contents of a scene where the nodes represent the detected objects, and edges linking them represent the inter-node relationships. An example scene graph is shown in Fig. 6, while a typical scene graph generation process is illustrated in Fig. 7.

The scene graph generation process consists of several steps [112]. First, given an image, an object detection module extracts object region proposals and visual features corresponding to these regions. Generally, a Faster-RCNN [67] is used for the backbone of object detection. The extracted features are used to identify object categories and their attributes. Identified objects along with their extracted features are used as nodes in the initial graph, e.g., a fully connected dense graph. Then, the extracted features along the nodes and edges are iteratively refined, and a final graph is inferred according to these refined features. Some recent papers [32, 31, 113] also consider joint optimization of object detection and relationship recognition parts. Specifically, Factorizable-Net is proposed in [33] where an RPN is used to extract object proposals and proposed objects are paired to obtain a fully-connected initial coarse graph. In [34], 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 [114] refines the features on the graph. On the other hand, VCTree model [115] constructs a dynamic tree from a scoring matrix where visual context is encoded into the tree structure. An illustration of the VCTree model is given in Fig. 8. Furthermore, some recent works [32, 31] use recurrent architectures for graph inference. Particularly, in [32] a feature refining
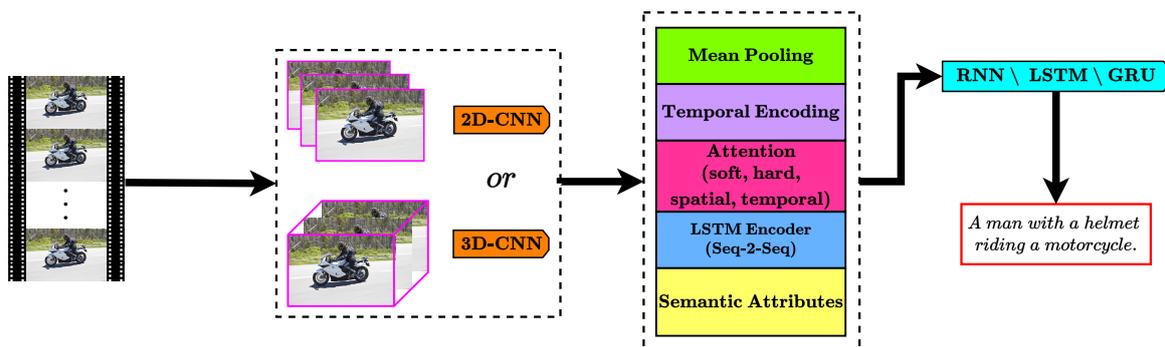


Figure 5: Video captioning pipeline. Often, frames are processed by 2D or 3D convolutions to extract features. Then, extracted features are processed further in the temporal domain by methods like attention, temporal encoding, LSTM, etc. Finally, a caption is generated by feeding recurrent architectures with temporally assessed features.
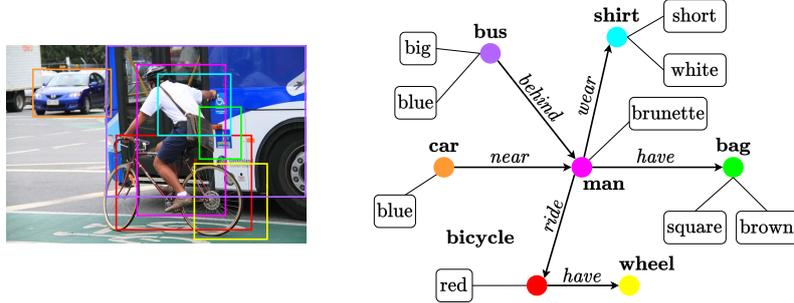
Figure 6: A scene graph example. Detected objects are represented as objects with circular colored nodes. Relationships among object nodes are shown with directed edges. Object attributes are denoted with rectangular uncolored nodes and they are connected to object nodes with undirected edges.
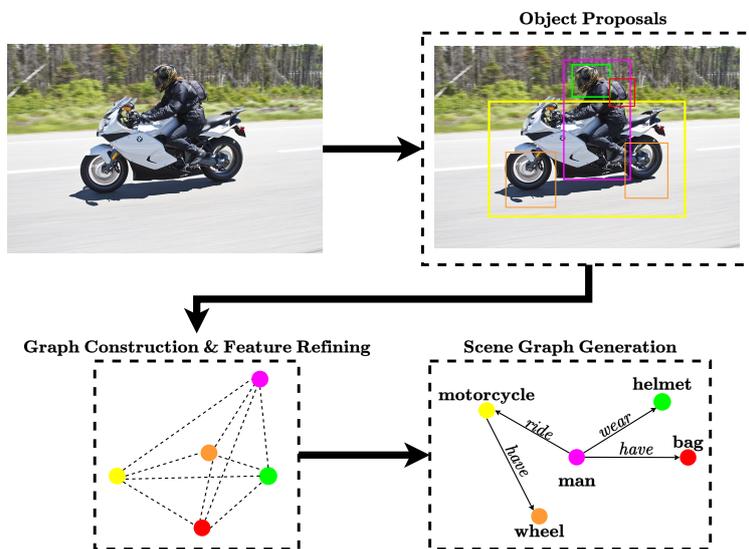


Figure 7: Scene graph generation pipeline. As a first step, object proposal regions are generated via RPN [67]. The proposals are used to construct a fully connected coarse graph. The initial graph and its features are refined iteratively. Refined features are used to infer node types and relationship types in the final graph.

module consisting of edge and node Gated Recurrent Units (GRU), and in [116], stacked bi-LSTMs are used.

There are approaches to incorporate external prior knowledge for scene graph generation tasks as well. In [118], natural language priors are incorporated and visual relationships and textual relationships are jointly learned and aligned. A linguistic knowledge distillation framework is proposed in [119] where statistics obtained from external texts are used to regularize visual models. In [120], knowledge-graphs are employed as prior information where the generated scene and knowledge graphs are abridged and iteratively refined.

Scene graphs can also be extracted from video signals. In [121], 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 [122] where scene-centric and view-centric graphs are hierarchically generated. For more details on scene graphs, we refer the readers to the pertinent survey papers [112, 123].
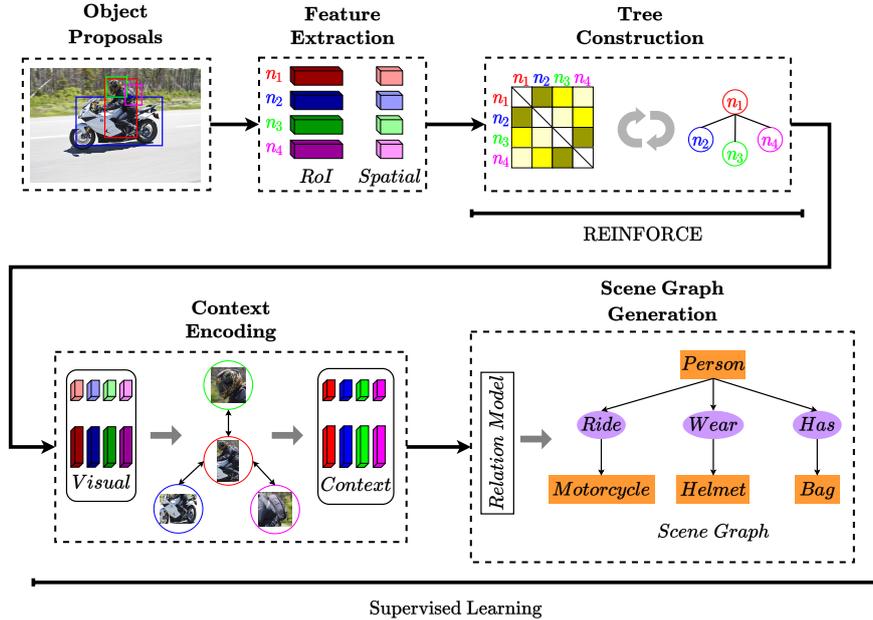
Figure 8: VCTree model proposed in [115]. Visual features are extracted from object proposals. Extracted features are used to compute a scoring matrix. Based on the score matrix, a dynamic tree structure is constructed using REINFORCE algorithm [117]. Finally, visual features are encoded into context features and a scene graph is generated via supervised learning.

### 3.1.4. Automatic Speech Recognition

Speech is one the most natural sources of semantic information. The most popular application of semantic transformation on speech signals is the conversion of audio signals into NL texts via automatic speech recognition (ASR) systems [124]. Similar to the image/video captioning applications, the range of target semantic modality is formed by a predetermined portion of the desired NL. The extent of the covered portion of the target natural language is continuously expanding throughout the years from digits [125] to tens of thousands of words [126].

Even though the requirements of an ASR system may vary for different applications (e.g., speaker dependency, vocabulary size, and utterance-awareness), most of the approaches conform with the process depicted in Fig. 9, as described in [124]. In Fig. 9, the input audio signal is first preprocessed via framing, filtering, Discrete Fourier Transform (DFT), denoising, or normalization, etc. Then, using the preprocessed version of the audio signal, some spectral or temporal features can be extracted via Mel-frequency Cepstral Coefficients (MFCC) [126–128] or Discrete Wavelet Transform (DWT) [129–131]. The extracted features are passed through a prediction module that employs Hidden Markov Models (HMMs) [132, 133], SVMs [134–136], RNNs [137–139], or CNNs [140–143], among others, to obtain the text equivalence in the desired language restricted by a predefined vocabulary and grammar rules. More details about ASR can be found in the pertinent survey papers [124, 144].

### 3.2. State-of-the-Art Semantic and Goal-Oriented Communications

With the tremendous advances in machine learning, particularly, on DNNs, there has been a renewed interest in semantic communications in recent years. As a recent example, joint learning and communication problem has been studied for multi-agent reinforcement learning (MARL) framework in [145], where the authors develop a systematic structure for collaborative MARL over noisy channels while taking into account the learning and the communication issues. In particular, a multi-agent partially observable Markov decision process (MA-POMDP) is studied in which agents can communicate with one another over a noisy channel to increase their shared long-term average reward. The agents not only learn to collaborate through a noisy link but also figure out how to communicate *effectively*, which is achieved by introducing the channel as a

Figure 9: Automatic speech recognition pipeline. The semantic features of a preprocessed speech audio signal are extracted and fed into a predictor for text generation. The prediction module is designed for a target language model.

part of the environment dynamics. Thus, the authors of [145] illustrate the increased performance of joint learning and communication framework compared to separate treatment of communications and the MARL problem. The *effectiveness problem* addressed partly by [145] was previously discussed in Section 2.

Another interesting study on semantic communications is performed in [146], where the authors focus on joint semantic channel coding by designing a deep-learning enabled semantic communication (DeepSC) system for text transmission. Unlike the traditional approaches, the objective of this study is to minimize the semantic errors in the transmitted text, as opposed to bit errors, while maximizing the system capacity. The performance metrics introduced by [146] rigorously formulate the semantic information at a sentence-level, and the corresponding semantic errors. The proposed *semantic level* communication architecture is shown to be more robust compared to the traditional *technical level* approaches, especially in the low signal-to-noise ratio (SNR) regime. In [51], a lite version of this distributed semantic communication system for text transmission is studied, where the authors consider an affordable IoT network by pruning the model redundancies and reducing the model size.

Transmission of speech data for semantic communications is considered in [147]. The authors propose a deep learning-enabled semantic communication system for speech signals (DeepSC-S). DeepSC-S extracts essential semantic information from speech audio signals by using squeeze-and-excitation (SE) networks [148]. Then, selective weighting is applied to emphasize the more essential information during DNN training. A joint design of speech encoder/decoder and the channel encoder/decoder is performed to learn and extract the speech features that mitigate the channel distortion.

In [22], the authors propose using semantic and goal-oriented communications in the next generation (6G) networks, which may significantly improve the effectiveness and sustainability of the system. The authors argue that instead of focusing on the bit-level recovery of the transmitted data, one may identify the relevant meaning, which can increase the effectiveness of the network without increasing the usage of communication resources such as bandwidth and energy. In [48], authors point out possible advantages of a goal-oriented joint consideration of information generation, transmission, and usage resulting in semantic-empowered communications, which is defined by transmitting the most informative data samples, enabling the end-user to attain the most current and essential information for its goal. Especially for future massive and intelligent systems, semantics-empowered communications are expected to improve the network resource utilization. Other studies on semantic communications are presented in [149–152].

Inspired by the previous work on semantic information, signal processing, and communications, we introduce a formal definition of the proposed goal-oriented semantic language and the semantic signal processing framework in the next section.

## 4. Proposed Goal-Oriented Semantic Language and Signal Processing Framework

In this section, we propose and rigorously define a graph-based semantic language to represent any type of signal in terms of a well-organized and easy-to-parse structure. We then describe a typical semantic signal processing framework and its basic building blocks in detail.

We propose a directed bipartite graph, shown in Fig. 10, as a semantic representation of a signal that consists of components with semantic connotations. In this representation, the signal components are first

11

classified into a predefined set of classes (represented as circles in Fig. 10). Then, their states (e.g., $c_1 \rightarrow p_0$) and relationship with other components (e.g., $c_1 \rightarrow p_1 \rightarrow c_2$) are labelled with appropriate predicates, again chosen from a predefined set (shown as squares in Fig. 10). As illustrated in the bipartite graph structure in Fig. 10, the signal components are depicted as nodes with labels $c_i$, and their semantic states/relationships are shown as nodes with labels $p_i$. Note that such a representation of raw sensor data might require a preprocessing stage that can involve transformation to an appropriate domain prior to the classification of signal components and identification of their semantic relationships. We further note that, unlike the general graph structures where extraction, pattern matching, etc., can be challenging, especially as the graphs get more complex, the proposed hierarchical bipartite structure allows for a complete but much simpler semantic description of a signal, and therefore, it is much simpler to generate and operate on.
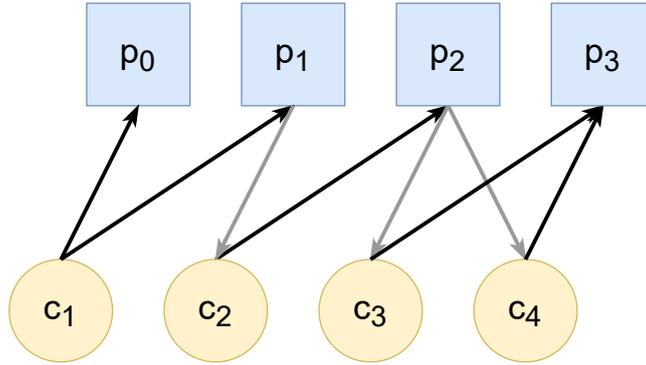


Figure 10: Proposed bipartite graph structure with signal components and predicates.

The proposed semantic description of a signal (e.g., a sensor output) includes a multi-graph description and corresponding attribute sets for each node in each graph. This structure of information in a signal enables easy parsing for goal-oriented storage, processing, and communications. In the following subsections, the multi-graph description, the companion attribute sets, and the goal-filtering process are presented in detail.

### 4.1. Multi-Graph Semantic Description

The node sets for components and predicates in the proposed bipartite graph structure are defined as follows:

$$C = \{c_1, c_2, \ldots, c_{N_c}\}, \tag{1}$$
$$P = \{p_0, p_1, p_2, \ldots, p_{N_p-1}\}, \tag{2}$$

where $N_c$ and $N_p$ are the numbers of component and predicate classes in the graph language, respectively. Note that a separate $p_0$ class is defined as the first element of the predicate classes. If a detected component has no other semantic connection to another component, then that component is connected to the predicate $p_0$ to form a bipartite graph. This is to ensure that a detected component never becomes an isolated node in a graph and always has a non-zero edge weight in the corresponding biadjacency matrix representation. Moreover, the predicate $p_0$ and the similarly defined predicates in the set $P$ can be considered as describing states of each component without describing any relation to the other components. Note that state describing predicates like $p_0$ have a single directed connection to a component, while relation predicates can have different polarities of connections depending on their definition (e.g., $c_i \rightarrow p_j \rightarrow c_k$ or $c_i \rightarrow p_j \leftarrow c_k$), as shown in Fig. 10.

Using the class definitions from (2), a *multi-graph class description* of a signal at time $t$ can be defined as a union of connected bipartite graphs as

$$\mathcal{S}_t = \{S_{t,1}, S_{t,2}, \ldots, S_{t,N}\}, \tag{3}$$

where $N$ is the number of *class atomic graphs* representing the disjoint bipartite graphs in the data, and

$$S_{t,i} = (C, P, E_{t,i}), \tag{4}$$

is a connected directed bipartite graph with node sets $(C, P)$ and edges (connections) at time $t$, $E_{t,i}$. Due to the bipartite nature of the proposed graph formalism, only connections from the component classes to the predicate classes and vice versa are allowed. Therefore, each connected graph $S_{t,i}$ in (3) can also be represented by the following biadjacency matrix:

$$\boldsymbol{S}_{t,i} = \begin{bmatrix} \mathbf{0} & \boldsymbol{S}_{t,i}^{CP} \\ \boldsymbol{S}_{t,i}^{PC} & \mathbf{0} \end{bmatrix}, \tag{5}$$

where $\boldsymbol{S}_{t,i}^{CP}$ and $\boldsymbol{S}_{t,i}^{PC}$ represent the connections from the components to the predicates and vice versa, respectively. In other words,

$$\boldsymbol{S}_{t,i}^{CP}[m,n] = \begin{cases} 1, & \text{if } \exists \text{ connection from } c_m \text{ to } p_n, \\ 0, & \text{otherwise}, \end{cases} \tag{6}$$

$$\boldsymbol{S}_{t,i}^{PC}[m,n] = \begin{cases} 1, & \text{if } \exists \text{ connection from } p_m \text{ to } c_n, \\ 0, & \text{otherwise}. \end{cases} \tag{7}$$

Note that in the multi-graph definition given in (3),(4) the nodes $C$ and $P$ are class definitions that may have multiple detections of their particular members in a signal. This class-only representation is a higher level of abstraction of the signal, and it is useful when the goals are defined similarly over classes and not over individual instances of components or predicates. For the case when the goals are defined over individual members, we introduce a lower level of abstraction again in the form of bipartite graphs to represent only the detected instances from classes $C$ and $P$. For the $i$-th atomic class graph $S_{t,i}$ at time $t$, we define the *detected component* and *detected predicate* sets as follows:

$$C_{t,i}^{D} = \{(c_j, k) \mid k \in [1, N_{c_j}^{t,i}] \text{ and } \exists m \; \boldsymbol{S}_{t,i}^{CP}[j,m] = 1 \text{ or } \boldsymbol{S}_{t,i}^{PC}[m,j] = 1\}, \tag{8}$$

$$P_{t,i}^{D} = \{(p_j, k) \mid k \in [1, N_{p_j}^{t,i}] \text{ and } \exists m \; \boldsymbol{S}_{t,i}^{CP}[m,j] = 1 \text{ or } \boldsymbol{S}_{t,i}^{PC}[j,m] = 1\}, \tag{9}$$

where $N_{c_j}^{t,i}$ and $N_{p_j}^{t,i}$ are the number of detected instances of component $c_j$ and predicate $p_j$ classes in the atomic graph $S_{t,i}$, and $k$ is the unique detection index for each class instance. Using the *detected* component and predicate sets, we can now define a *multi-graph instance representation* as

$$\mathcal{D}_t = \{D_{t,1}, D_{t,2}, \ldots, D_{t,N}\}, \tag{10}$$

$$D_{t,i} = (C_{t,i}^{D}, P_{t,i}^{D}, E_{t,i}^{D}), \tag{11}$$

where $D_{t,i}$ are the *instance atomic graphs*. The corresponding connection lists $E_{t,i}^{D}$ are defined as triplets indicating a connection from $(c_a, k_a)$ to $(p_b, k_b)$ to $(c_c, k_c)$ as

$$E_{t,i}^{D} = \{((c_a, k_a), (p_b, k_b), (c_c, k_c))\}. \tag{12}$$

Note that in (12), the last entry of each triplet can be an empty element, as there can be a single connection from a component to a predicate, but not vice versa. This is in line with the definition of the $p_0$ predicate in (2). The biadjacency matrix representations $\boldsymbol{D}_{t,i}, \boldsymbol{D}_{t,i}^{CP}, \boldsymbol{D}_{t,i}^{PC}$ for $D_{t,i}$ can be defined similarly to (5)–(7).

As noted before, the multi-graph sets $\mathcal{S}_t$ and $\mathcal{D}_t$ provide a hierarchical bipartite representation that is much simpler to generate and operate on than a general single graph structure. In typical machine-type applications, the number of components and predicates scale according to the complexity and the processing capabilities of the sensors or devices. For example, a simple IoT sensor may not have more than a few components and predicates, which—combined with the multi-graph representations—makes pattern matching and goal-filtering operations extremely simple implementations of the available methods in the literature [59–62].

## 4.2. Attribute Sets in Multi-Graph Semantic Descriptions

The class definitions and detected class instances do not necessarily constitute a complete semantic description of a signal. As such, to fully capture additional properties and features, we use associated attribute sets for each component or predicate node that is present in each atomic graph $D_{t,i}$. For each $D_{t,i}$ in the multi-graph description of the signal, we define an attribute superset $A_{t,i}$ as

$$A_{t,i} = \{\Theta_{t,i}(n_j, k) \mid (n_j, k) \in C_{t,i}^D \cup P_{t,i}^D\}, \tag{13}$$

where $(n_j, k)$ corresponds to a node (component or predicate) and instance index pair. The operator $\Theta_{t,i}(.)$ is defined as the multi-level attribute set of its argument nodes as

$$\Theta_{t,i}(n_j, k) = \{\boldsymbol{\theta}_{t,i}^{(1)}(n_j, k), \boldsymbol{\theta}_{t,i}^{(2)}(n_j, k), \ldots, \boldsymbol{\theta}_{t,i}^{(L_{n_j})}(n_j, k)\}, \tag{14}$$

where $L_{n_j}$ is defined as the number of levels in attributes for node-$j$, which may be defined separately for each type of component or predicate, or kept the same for all nodes for simplicity. The introduction of levels of attributes enables a hierarchical organization of attributes from a simpler to a more complex representation, which in turn makes goal-oriented filtering easier and more intuitive. The exact definition of (14) depends on the application, the nature of internal/external goals, and the type of signal processing to be performed on the obtained semantic information. Concrete examples of how attribute sets can be defined are given via several use-cases in Section 5.

The companion attribute set for the object multi-graph description $\mathcal{D}_t$ is defined as

$$\mathcal{A}_t = \{A_{t,1}, A_{t,2}, \ldots, A_{t,N}\}. \tag{15}$$

Finally, the full semantic description of the sensor data at time $t$ is defined as

$$Y_t = (\mathcal{S}_t, \mathcal{D}_t, \mathcal{A}_t), \tag{16}$$

where $\mathcal{S}_t$, $\mathcal{D}_t$, and $\mathcal{A}_t$ correspond to the multi-graph class description, multi-graph instance description, and the corresponding attribute sets, respectively.

## 4.3. Semantic Goals and Goal-Oriented Filtering

We now define a goal operator using the proposed semantic framework to parse and filter the extracted semantic information from a signal. Note that beyond this point, the time index subscript $t$ is dropped for simplicity. A goal-filtering operator $\mathcal{G}$ acting on semantic description of the sensor data $Y$ is defined as

$$\mathcal{G}(G^S, G^D, \boldsymbol{l}_g, \boldsymbol{l}_a)\{Y\} = \widetilde{Y}, \tag{17}$$

$$\widetilde{Y} = (\widetilde{\mathcal{S}}, \widetilde{\mathcal{D}}, \widetilde{\mathcal{A}}), \tag{18}$$

where the terms inside the parentheses in (17) are the arguments of the goal operator $\mathcal{G}$. $\widetilde{Y}$ is the filtered semantic description, and $(\widetilde{\mathcal{S}}, \widetilde{\mathcal{D}}, \widetilde{\mathcal{A}})$ correspond to the goal-filtered versions of $(\mathcal{S}, \mathcal{D}, \mathcal{A})$, respectively. $G^S = (C, P, E^S)$ and $G^D = (C^D, P^D, E^D)$ are the bipartite graphs that correspond to the queried patterns in the class-level graph $\mathcal{S}$ and the instance-level graph $\mathcal{D}$, respectively. The last two arguments are the graph

complexity vector $\boldsymbol{l}_g$ and attribute complexity vector $\boldsymbol{l}_a$ that define the desired level of complexity for the semantic output in terms of the multi-graph representation and the attribute sets, respectively.

The distinction between class-level and instance-level queries for $\mathcal{G}$ enables the goals to be defined at different levels of abstraction; namely, *global goals* with a typical interest on class-level information, and *local goals* with a typical interest on detected instances and their relations. Hypothetically, the global goals can be defined first and change slowly over time, while local goals can be defined following the detection of a component of interest and change more rapidly in time.

The graph complexity vector $\boldsymbol{l}_g = \{l_g \in \mathbb{N}\}$ is a vector of natural numbers $l_g$, defining an $l_g$-hop neighborhood around components of each queried graph pattern in $\mathcal{G}$. This allows goals to search for a specific pattern with $G^S$ or $G^D$, while being interested in other relationships among the components in those specific patterns without explicitly defining them.

Using $G^S$, $G^D$, and $\boldsymbol{l}_g$, the output semantic graphs $\widetilde{\mathcal{S}}$, $\widetilde{\mathcal{D}}$ can be defined as

$$\widetilde{\mathcal{S}} = \{\widetilde{S}_i \mid \widetilde{S}_i \subset S_i, G^S \subset \widetilde{S}_i \text{ and } \widetilde{S}_i \text{ includes } l_g\text{-hop neighborhood of } G^S \text{ within } S_i\}, \tag{19}$$

$$\widetilde{\mathcal{D}} = \{\widetilde{D}_i \mid \widetilde{D}_i \subset D_i, G^D \subset \widetilde{D}_i \text{ and } \widetilde{D}_i \text{ includes } l_g\text{-hop neighborhood of } G^D \text{ within } D_i\}. \tag{20}$$

Note that the specific pattern matching and $l$-hop neighborhood search algorithms employed to generate (19) and (20) can be chosen from many available alternatives in the literature [59–62]. However, we emphasize again that the proposed multi-graph structure allows for a simple and complete description of a signal, compared to a general single-graph structure. Therefore, the computational complexity of the pattern matching and search algorithms is relatively low and they scale according to the complexity of the sensor or device. An illustration of the goal filtering with a given graph complexity of $l_g = 1$ is given in Fig. 11. A simple class-level goal $(c_2 \to p_1)$ is applied to the instance representation in Fig. 11b, with $l_g = 1$. In Fig. 11c, the exact matching of the goal is shown with red edges, and the requested neighborhood around the matched pattern is shown inside the shaded parallelogram. Note that the graph complexity parameter is defined over neighboring components. Therefore, the output pattern includes components separated by at most $l_g$ predicates with the exact-matched pattern.



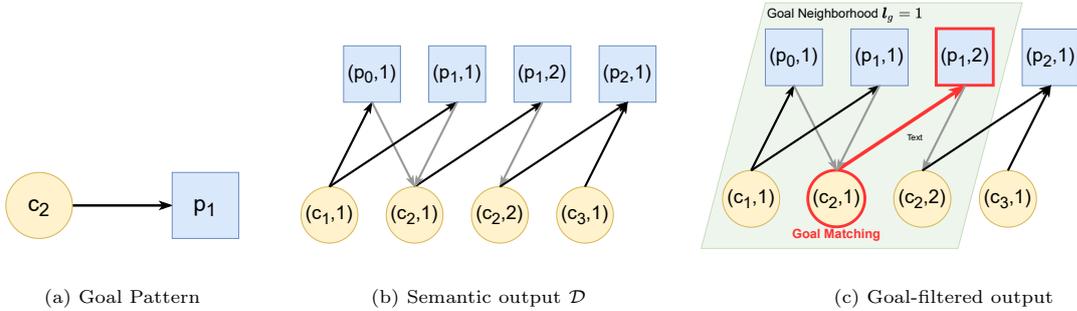| (a) Goal Pattern | (b) Semantic output $\mathcal{D}$ | (c) Goal-filtered output |

Figure 11: Goal filtering example. Note that $l_g = 1$ covers neighboring components that are connected to the matched pattern by a single predicate path.

After defining the filtered graph outputs $\widetilde{\mathcal{S}}$ and $\widetilde{\mathcal{D}}$, their corresponding attribute sets are filtered with the attribute complexity vector $\boldsymbol{l}_a = \{l_a \in \mathbb{N}\}$ that requests only the first $l_a$ levels within the hierarchy of attribute sets of each instance in the graphs as

$$\widetilde{\mathcal{A}} = \{\widetilde{\Theta}_i \in \widetilde{D}_i \mid \Theta_i = \{\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(L_{min})}\}\}, \tag{21}$$

where $L_{min} = min(l_a, L)$ is the minimum of the desired and available levels of complexity.
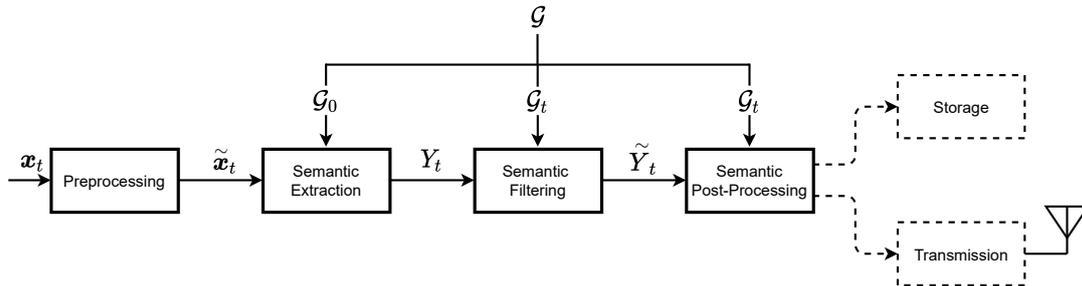
Figure 12: The proposed goal-oriented semantic signal processing framework.

## 4.4. Goal-Oriented Semantic Signal Processing Framework

Based on the goal-oriented semantic language framework defined above, the proposed semantic signal processing framework in its generic form is illustrated in Fig. 12. In the following, we give brief descriptions for each conceptual block in Fig. 12, before discussing its potential implementations at the goal and the hardware levels.

The proposed goal-oriented semantic signal processing framework consists of the conceptual blocks of preprocessing, semantic extraction, semantic filtering, semantic post-processing, and storage or transmission, depending on the application. In the *preprocessing* block, following a prefiltering for noise and interference reduction, the input signal $\boldsymbol{x}_t$ is transformed to an appropriate domain where detection and classification of components and predicates can be performed efficiently. For one-dimensional (1-D) input signals, decomposition of the data for graphical representation may not be as self-evident as in the case for higher-dimensional input data. In such cases, commonly used domain transformation approaches such as the time-frequency and the time-scale domain processing [153–155], can be used to segment and detect distinct signal components via pattern recognition and computer vision techniques [156]. For image and video applications, transformation to another domain may not be necessary, as component detection and classification in these domains can be efficiently implemented on the input signals directly.

The *semantic extraction* block is where the multi-graph description and corresponding attribute sets are generated from the input data. In the case of semantic representations of 2-D signals such as images or 3-D data such as video streams of consecutive slices of MRI data, components or objects can be first identified by image segmentation. For this purpose, scene graph extractors can be used [31–36]. The predicates can be identified based on the spatial relationships among the objects in an image, or their spatio-temporal relationships in a video stream. Note that a global goal $\mathcal{G}_0$ can be incorporated at this point. If a constant or slowly changing goal is defined either internally, or it is received from an external agent, the semantic extraction algorithm can adapt accordingly. For example, if the global goal is only related to a small subset of potential components that can be extracted, the algorithm efficiency can be increased by limiting the output classes of the classification algorithm implemented as a DNN.

In the next block, *semantic filtering* operations are performed. Once the scene descriptions $Y_t$ are generated, local and time-varying goals $\mathcal{G}_t$ can be applied to filter *interesting* semantic information within $Y_t$. In a typical application, local goals $\mathcal{G}_t$ are either generated internally via decision-making algorithms or are received from external agents. This block typically employs graph signal processing algorithms to match and extract goal queries from the graph-based representation of the signal. Note that the semantic output $Y_t$, generated by the proposed language yields relatively simple graphs due to their bipartite and goal-oriented definitions; hence, the computational complexity of any graph operation is considered to be more tractable compared to general graph structures.

The *semantic post-processing* block is a conceptual block where further processing of the extracted and goal-filtered semantic data can be performed. Note that at this point, the semantic data are strictly organized and only include goal-filtered components and their corresponding attributes. These attributes allow the implementation of a wide spectrum of processing tasks. For instance, since the location information of the components is a probable part of the attribute sets, it can allow for tracking of components or groups of

16

components by using standard tracking techniques. Furthermore, the detailed attributes of components such as wavelet representation of its time-scale characterization facilitate the application of more specific signal processing algorithms on the components of interest such as further denoising or statistical modeling. The semantic post-processing block can also incorporate local goals $\mathcal{G}_t$, to schedule transmission and storage operations according to the level of desired semantic information.

Finally, the *storage* and the *transmission* are tentative blocks that can be added to the framework depending on the application. In either case, an encoding scheme is utilized to compress the goal-filtered graphs and attributes. Then, the compressed semantic information can either be stored internally or passed forward to a transmission block where further channel coding operations can be performed prior to transmission.

The semantic signal processing framework described above can be adopted in the next generation of devices and communication networks. Although the types and implementation strategies of goals are previously discussed, the generation and dissemination of semantic goals warrant further discussion. Typically, global goals $\mathcal{G}_0$ for a semantic device can be defined either by human operators at the language level or by manufacturers at a hardware level. Alternatively, in a hierarchical network such as a Smart City application [157, 158], global goals can be defined at the highest level, and then they can be disseminated by mapping the global goals to the low-level device language by intermediate devices such as base stations. On the other hand, the generation of rapidly varying local goals $\mathcal{G}_0$ requires decision-making algorithms to be implemented at intermediate levels to optimize throughput and goal implementation. Possible approaches to the goal generation and dissemination problems address the so-called *effectiveness problem* described in Section 2, and constitute a future research direction for goal-oriented semantic signal processing.

The adoption of the proposed semantic signal processing framework in Fig. 12 will also have repercussions at a hardware level. The preprocessing and the semantic extraction stages can be implemented as part of a sensor subsystem designed to acquire signals of the desired modality such as audio or video signals. Such a sensor subsystem must provide its output in the proposed semantic structure so that the rest of the semantic processing chain can be implemented on its output. Therefore, the proposed semantic signal processing framework may form a basis for sensor standardization efforts for future applications of semantic signal processing and communication systems.

In the following sections, detailed examples of generating and processing semantic information within the proposed framework as well as possible encoding/transmission strategies are discussed in detail.

## 5. Applications of Semantic Signal Extraction with the Proposed Framework

In this section, we provide concrete examples of the proposed goal-oriented semantic signal processing framework described above. First, we present a real-time computer vision application for a video stream input. Next, in contrast to the higher dimensional nature of video signals, the proposed language is implemented on a generic 1-D sensor output on which we also implement a novel and efficient non-uniform sampling strategy. We then provide a brief discussion on the potential applications of the proposed framework on heterogeneous multi-sensor networks. It is important to note that these case studies are provided only as a proof-concept of the generality of the proposed framework, which has a great potential to be customized for many different applications outside of what is presented here.

### 5.1. Real-Time Semantic Processing of Video-Stream Data

As a first application scenario, we demonstrate the use of the proposed semantic signal processing framework on real-time video signals via the architecture shown in Fig. 13. An input sequence of frames denoted as $F_n$, with $n$ as the discrete time index, is fed into a detection and classification block, and the detected entities in each frame are filtered out via a global goal $\mathcal{G}_{object}$. The remaining *globally-interesting* objects and extracted predicates among them are tracked across frames using their position and velocity information provided by the object tracking block. Based on the tracks, a *multi-graph instance representation* $\mathcal{D}_n = \{D_1, \ldots, D_{M_n}\}$ composed of $M_n$ disconnected subgraphs and a corresponding attribute set $\mathcal{A}_n = \{A_1, \ldots, A_{M_n}\}$ are generated. The detailed semantic representation for each frame is updated in accordance with the tracking unit's updates. In this specific application, the *multi-graph class representation*
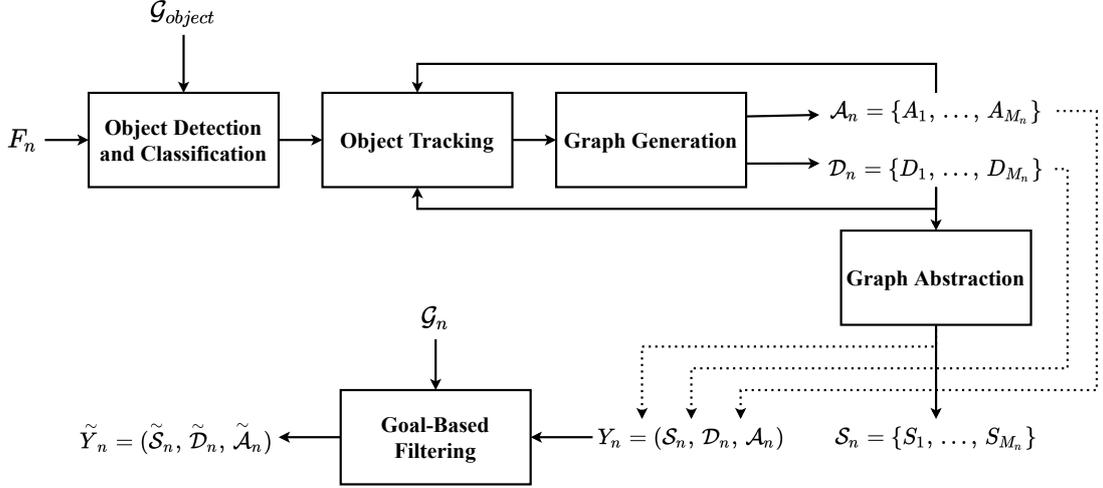
Figure 13: Semantic processing architecture for real-time video signals. Objects are detected and classified at each frame $F_n$ by YOLOv4-CSP [18]. For temporal extension, detected objects are tracked across the frames using DeepSORT [159]. $\mathcal{S}_n$ is obtained by performing a summarization of $\mathcal{D}_n$, i.e., simply drawing a crude version of $\mathcal{S}_n$ by only keeping the class information of the nodes. Finally, the complete semantic description $Y_n$ consists of the generated graph $\mathcal{D}_n$, its abstract version $\mathcal{S}_n$ and its attribute set $\mathcal{A}_n$. The resulting semantic description $Y_n$ is then filtered via local goals $\mathcal{G}_n$.

$\mathcal{S}_n$ is constructed via a post-processing step as a higher-level abstraction and summary of $\mathcal{D}_n$ in the Graph Abstraction block. Finally, the complete semantic description $Y_n = (\mathcal{S}_n, \mathcal{D}_n, \mathcal{A}_n)$ is generated. The local goals $\mathcal{G}_n$ prune the semantic description $Y_n$ to generate the goal-oriented semantic description $\widetilde{Y}_n$. Note that, in this example, we are considering fast algorithms that are available for real-time implementation. Alternatively, for delay-tolerant (offline) applications, scene graph generator models [123] can be used directly in place of the first three blocks. A detailed discussion of each block and the algorithms employed are discussed in the following.

Taking into account the real-time processing requirements, we use a pretrained YOLOv4-CSP [18] model on COCO dataset [160] for our object detection module since it achieves greater accuracy compared to the other models [70, 161] under the latency requirements of a typical video signal with 24 frames-per-second. Consequently, our semantic graph language contains $N_C = 80$ semantic component classes, i.e., $C = \{c_1, \ldots, c_{80}\}$, corresponding to the object categories in the COCO dataset ($c_1$ : person, $c_2$ : bike, $c_{18}$ : dog, etc.). Omitting the time index $n$ for simplicity, YOLOv4-CSP provides each detected object in the form $o_i = (\boldsymbol{b}_i, \boldsymbol{y}_i)$, where $\boldsymbol{b}_i \in \mathbb{R}^4$ contains the position of the object according to its bounding box and $\boldsymbol{y}_i \in \mathbb{R}^{N_C}$ contains the category confidence scores. Here, we determine the classes of detected semantic instances via $\boldsymbol{y}_i$ and represent each object in the form $o_i \equiv (c_a, i)$. This representation simply indicates that the object $o_i$ is an instance of component class $c_a$ with a unique identifier index $i$ where

$$a = \arg \max_{c=1,\ldots,80} \boldsymbol{y}_i[c]. \tag{22}$$

To put it another way, $(c_a, i)$ is the pointer of object $o_i$ so that $C^D$ as defined in (8) stores the detected objects as

$$C^D = \left\{ (c_a, i) \equiv o_i \,\middle|\, c_a \in C \right\}. \tag{23}$$

An example output of object detection/classification module is shown in Fig. 14. The three detected objects belong to the semantic component classes $c_1$ : person, $c_{18}$ : dog, and $c_{19}$ : horse, i.e., the set of detected objects is simply $C^D = \{o_1 \equiv (c_1, 1), o_2 \equiv (c_{19}, 2), o_3 \equiv (c_{18}, 3)\}$.
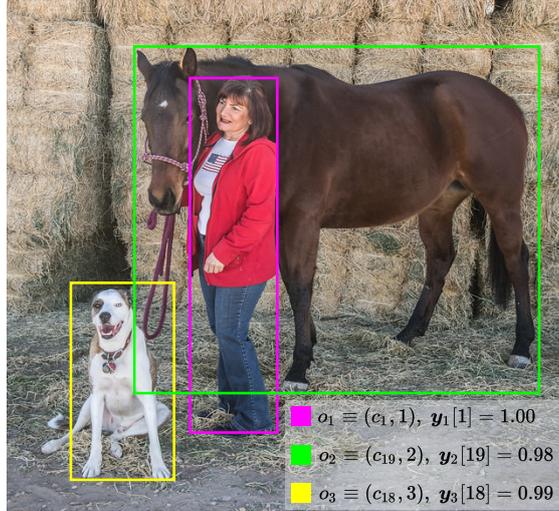
18

Figure 14: An example of detected objects by YOLOv4-CSP model. Predicted bounding boxes $\boldsymbol{b}_i$ are shown and maximum element of class confidence scores $\max\{\boldsymbol{y}_i[c]\}$ are listed.

In this example, we define our predicate set as $P = \{p_0, p_1, p_2\}$ where $p_0$ : null, $p_1$ : moving-together, and $p_2$ : conjunct. The *null* predicate $p_0$ is connected to isolated objects by default (see the discussion on avoiding non-zero adjacency matrices in Section 4). The predicate $p_1$ : *moving-together* is placed between adjacent objects moving towards to the same direction with similar velocities, whereas $p_2$ : *conjunct* predicate is placed between objects which have significantly overlapping bounding boxes and possibly possessive relationships, e.g., a person with their bike or bag. Note that all predicates are defined in such a way that all the edges start from the detected components and end at the predicates. This is due to the symmetric nature of the predicates in this specific application. Moreover, it should be noted that these predicates are only defined due to the simplicity of their inference, and the predicate set can be enriched further by extracting other relationships among the detected instances. For instance, if a *person* and a *bicycle* are moving towards the same direction with a significant overlap between their bounding boxes, then the predicate between these objects can be set to *riding*. With a similar reasoning for objects *person* and *car*, a *driving* predicate can be generated. Analogous to the objects, we denote the predicates as $e_k \equiv (p_w, k)$ with $k$ being a unique identifier index. Predicates $e_k$ are placed between some objects $o_i \equiv (c_a, i)$ and $o_j \equiv (c_b, j)$, where the set $P^D$ as defined in (9) stores the detected predicates as

$$P^D = \left\{ (p_w, k) \equiv e_k \,\middle|\, p_w \in P \right\}. \tag{24}$$

To identify the temporal and spatial changes within frames we perform tracking on the detected objects. Just like the object detection module, we require a tracker with a fast inference capability. Therefore, we use the DeepSORT [159] algorithm for tracking individual objects in the video stream. Specifically, DeepSORT utilizes a Kalman filter [162] to recursively predict future positions of the objects. The Kalman filter models the state transitions using a linear constant velocity model. We denote the state vector of each object $o_i \equiv (c_a, i)$ with $\boldsymbol{m}_i$, which contains the parameters of its bounding box $\boldsymbol{b}_i$ and its respective velocities as

$$\boldsymbol{m}_i = \left[ x_i^c, y_i^c, w_i, h_i, \dot{x}_i^c, \dot{y}_i^c, \dot{w}_i, \dot{h}_i \right], \tag{25}$$

where $x_i^c, y_i^c, w_i, h_i$ are the center coordinates, the width, and the height of the bounding box $\boldsymbol{b}_i$, respectively. To improve the tracking performance under challenging scenarios such as occlusion, non-stationary camera and multiple viewpoints, DeepSORT further utilizes deep cosine association metric [159, 163]. Particularly, each detected object $o_i \equiv (c_a, i)$ is cropped from the original frame based on the bounding boxes $\boldsymbol{b}_i$ and
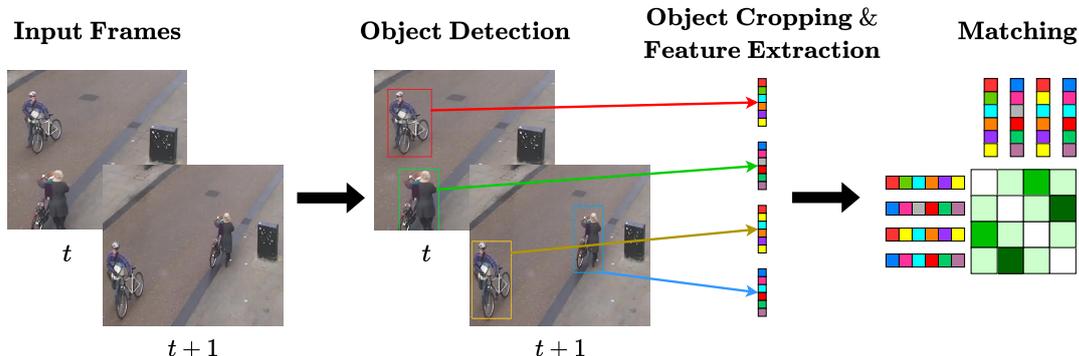
19

Figure 15: DeepSORT [159] compares feature vectors of detected objects and tracks them for track assignment. Detected objects are cropped and fed into a CNN for feature extraction. Features are utilized as additional information for track/detection matching.
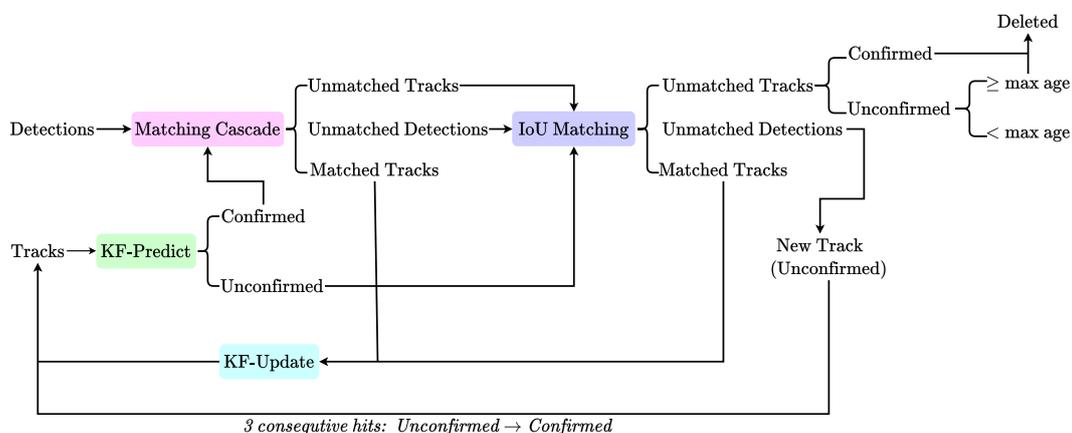


Figure 16: The DeepSORT tracking algorithm. At each frame, states of tracked objects are updated via a Kalman filter. Matching Cascade module matches detections with existing tracks using both extracted features and predictions of the Kalman filter. Specifically, detections and tracks are matched using Hungarian assignment [164]. The cost matrix of the Hungarian algorithm consists of the smallest cosine distances between extracted feature vectors and Mahalanobis distances between bounding boxes (see [159] for details).

passed through a CNN to be represented as unit-norm vectors $\boldsymbol{r}_i \in \mathbb{R}^{128}$, $\|\boldsymbol{r}_i\|_2 = 1$, as illustrated in Figs. 15 and 16. We refer to $\boldsymbol{r}_i$'s as the feature vectors. In fact, for each tracked object $o_i \equiv (c_a, i)$, we store multiple feature vectors across frames $\left\{ \boldsymbol{r}_i^{(l)} \right\}_{l=1}^{T_i}$ where $\boldsymbol{r}_i^{(l)}$ stands for the $l$-th feature vector of the tracked object $o_i \equiv (c_a, i)$, and $T_i$ is number of times the object $o_i$ is observed.

We employ the state $\boldsymbol{m}_i$, feature vectors $\left\{ \boldsymbol{r}_i^{(l)} \right\}_{l=1}^{T_i}$ and the cropped RGB image $I_i \in \mathbb{R}^{w_i \times h_i \times 3}$ as the multi-level attributes of object $o_i \equiv (c_a, i)$ with $L = 3$ number of levels as defined in (14), i.e.,

$$\Theta(c_a, i) = \left\{ \boldsymbol{\theta}^{(1)}(c_a, i), \boldsymbol{\theta}^{(2)}(c_a, i), \boldsymbol{\theta}^{(3)}(c_a, i) \right\} = \left\{ \boldsymbol{m}_i, \{\boldsymbol{r}_i^{(l)}\}_{l=1}^{T_i}, I_i \right\}. \tag{26}$$

Note that, it is possible that different categories can have a different number of attribute levels $L_a$; however, for our experiments, we set $L_a = 3$, $\forall a$, for simplicity. It is important to note that, even though in our experiments, $\boldsymbol{\theta}^{(3)}(c_a, i)$ is set to the raw data $I_i$ itself, it can also be chosen as any multi-scale representation of $I_i$ such as its wavelet decomposition or discrete cosine transform. As an example, attributes for a detected *person* object is illustrated in Fig. 17.

(a) Frame under process, generated instance graph, and corresponding attribute set for Person-40.

(b) Lowest level attribute set corresponding to the state vector $\boldsymbol{m}$ of Person-40. Units are in pixels for position and pixels-per-frame for velocities.

(c) Second and third level attributes of Person-40. The second level includes a set of collected feature vectors across 48 frames with $T_{40} = 48$. The third level includes the cropped RGB image $I_{40}$, although alternative encoded versions of $I_{40}$ can also be used.
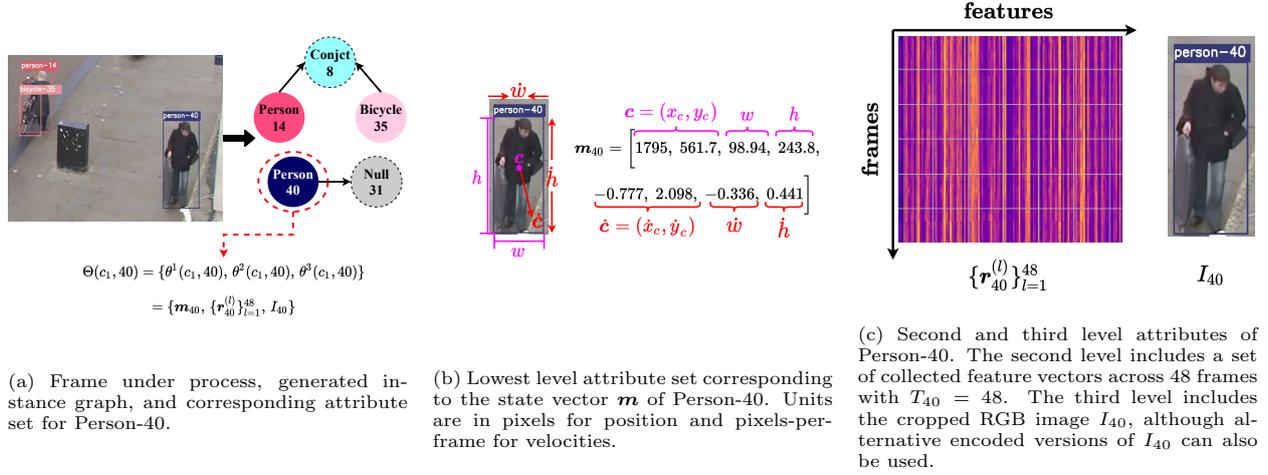
Figure 17: Illustration of attributes for a *person* component instance.

The exploitation of object positions and velocities enables us to determine types of predicates between different object pairs. We identify the type of a predicate between objects $o_i \equiv (c_a, i)$ and $o_j \equiv (c_b, j)$ by using their states $\boldsymbol{m}_i$ and $\boldsymbol{m}_j$. The following three metrics are defined for predicate detection

$$d^{(1)}(i,j) = \frac{\left\| (x_i^c, y_i^c) - (x_j^c, y_j^c) \right\|_2^2}{\sqrt{w_i \cdot h_i \cdot w_j \cdot h_j}}, \tag{27}$$

$$d^{(2)}(i,j) = \frac{(\dot{x}_i^c, \dot{y}_i^c)^\top (\dot{x}_j^c, \dot{y}_j^c)}{\left\| (\dot{x}_i^c, \dot{y}_i^c) \right\|_2 \left\| (\dot{x}_j^c, \dot{y}_j^c) \right\|_2} = cos\left( (\dot{x}_i^c, \dot{y}_i^c), (\dot{x}_j^c, \dot{y}_j^c) \right), \tag{28}$$

$$d^{(3)}(i,j) = IoU\left( \boldsymbol{b}_i, \boldsymbol{b}_j \right). \tag{29}$$

Here, $d^{(1)}(i,j)$ corresponds to scaled Euclidean (2-norm) distance between center coordinates of the positions, $d^{(2)}(i,j)$ is the cosine similarity between the velocities, while $d^{(3)}(i,j)$ is the intersection-over-union of bounding boxes of objects $o_i$ and $o_j$. Specifically, for a predicate $e_k \equiv (p_1, k)$, i.e., *moving-together* between objects $o_i \equiv (c_a, i)$ and $o_j \equiv (c_b, j)$, across multiple frames we require

$$\mathbb{1}\left( d^{(1)}(i,j) \leq z^{(1)} \right) \cdot \mathbb{1}\left( d^{(2)}(i,j) \geq z^{(2)} \right), \tag{30}$$

and similarly, for $e_k \equiv (p_2, k)$, i.e., *conjunct*, we require

$$\mathbb{1}\left( d^{(1)}(i,j) \leq z^{(1)} \right) \cdot \mathbb{1}\left( d^{(2)}(i,j) \geq z^{(2)} \right) \cdot \mathbb{1}\left( d^{(3)}(i,j) \geq z^{(3)} \right), \tag{31}$$

for predetermined thresholds $z^{(1)}$, $z^{(2)}$ and $z^{(3)}$, where $\mathbb{1}(\cdot)$ denotes the indicator function. Furthermore, even though it is not crucial for our application, we define attributes for a predicate $e_k \equiv (p_w, k)$, $p_w \in \mathcal{P} \setminus \{p_0\}$ as

$$\Theta(p_w, k) = \left\{ \boldsymbol{\theta}^{(1)}(p_w, k) \right\} = \left\{ \left[ d^{(1)}(i,j), d^{(2)}(i,j), d^{(3)}(i,j) \right] \right\}, \tag{32}$$

with $L_w = 1$, $\forall w$. Here, $\Theta(p_w, k)$ consists of a single-level attribute, i.e., a vector that contains metric evaluations between objects $o_i$ and $o_j$. Note that unlike the other predicates belonging to classes in $\mathcal{P} \setminus \{p_0\}$ which are connected to two objects $o_i \equiv (c_a, i)$ and $o_j \equiv (c_b, j)$, a null predicate $e_l \equiv (p_0, l)$ is only connected to $o_u \equiv (c_f, u)$, for an isolated object $o_u$. For null predicates, the attribute set is simply defined as the the empty set, i.e., $\Theta(p_0, l) = \varnothing$.

Given the set of detected objects $C^D$ and predicates $P^D$, we denote our detected connection set $E^D$ with triplets and pairs as

$$E^D = \left\{ (o_i, e_k, o_j) \,\middle|\, o_i, o_j \in C^D, e_k \in P^D, p_w \in P \setminus \{p_0\} \right\} \bigcup \left\{ (o_u, e_l) \,\middle|\, o_u \in C^D, e_l \equiv (p_0, l) \right\}. \quad (33)$$

In other words, the detected connection set $E^D$ simply consists of *(object, predicate, object)* triplets and *(unaccompanied object, null predicate)* pairs. Consequently, the graph $D$ is generated by the detected object, predicate, and edge sets as

$$D = \left( C^D, P^D, E^D \right). \quad (34)$$

We then define the companion attribute set $A$ of graph $D$ as

$$A = \left\{ \Theta(c_a, i) \,\middle|\, o_i \equiv (c_a, i) \in C^D \right\} \bigcup \left\{ \Theta(p_w, k) \,\middle|\, e_k \equiv (p_w, k) \in P^D \right\}, \quad (35)$$

which stores attributes for both object and predicate instances.

Furthermore, we generate the class-level representation of instance graph $D$ as defined in (4). The class description graph $S$ is defined as

$$S = (C^S, P^S, E^S), \quad (36)$$

where

$$C^S = \left\{ c_a \,\middle|\, o_i \equiv (c_a, i) \in C^D \right\}, \quad (37)$$

$$P^S = \left\{ p_w \,\middle|\, e_k \equiv (p_w, k) \in P^D \right\}, \quad (38)$$

$$E^S = \begin{aligned} &\left\{ (c_a, \, p_w, \, c_b) \,\middle|\, (o_i, e_k, o_j) \equiv ((c_a, i), \, (p_w, k) \, (c_b, j)) \in E^D \right\} \\ &\bigcup \left\{ (c_f, p_0) \,\middle|\, (o_u, e_l) \equiv ((c_f, u), \, (p_0, l)) \in E^D \right\}. \end{aligned} \quad (39)$$

To facilitate a better understanding of the high-level abstraction provided by (36), we provide an illustrative example in Fig. 18. The class level graph generation is a summarization of the components and predicates in the instance level graph, i.e., it is a surjective function mapping $D$ to $S$.

Finally, we define the semantic description $Y$ as a triplet consisting of the class-level graph, the instance-level graph, and attribute supersets as

$$Y = (\mathcal{S}, \mathcal{D}, \mathcal{A}). \quad (40)$$



(a) Frame under process with detections.

(b) Instance-level graph $D$. Object nodes are denoted with solid circles, predicates are denoted with dashed circles. Isolated objects are connected to the *null* predicate by default.

(c) Class-level graph $S$.

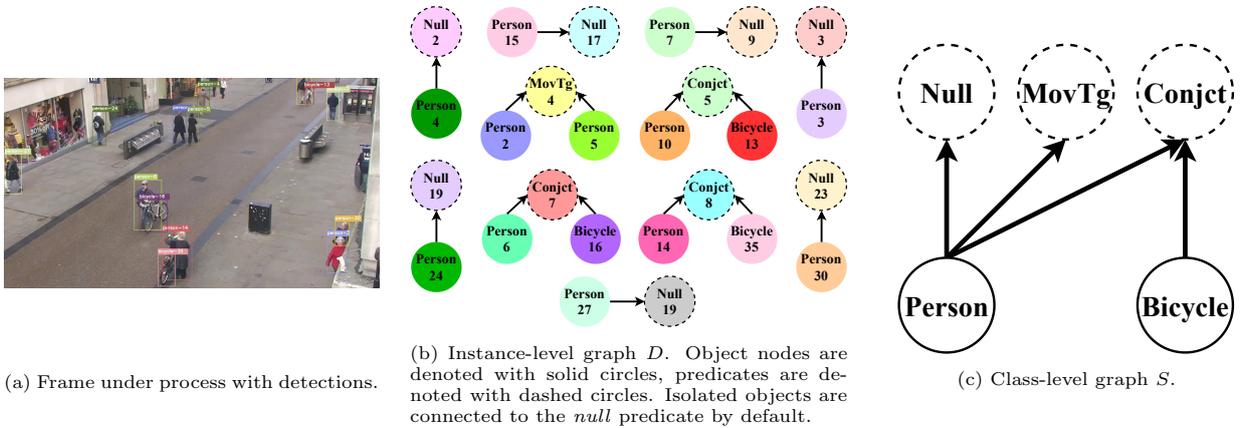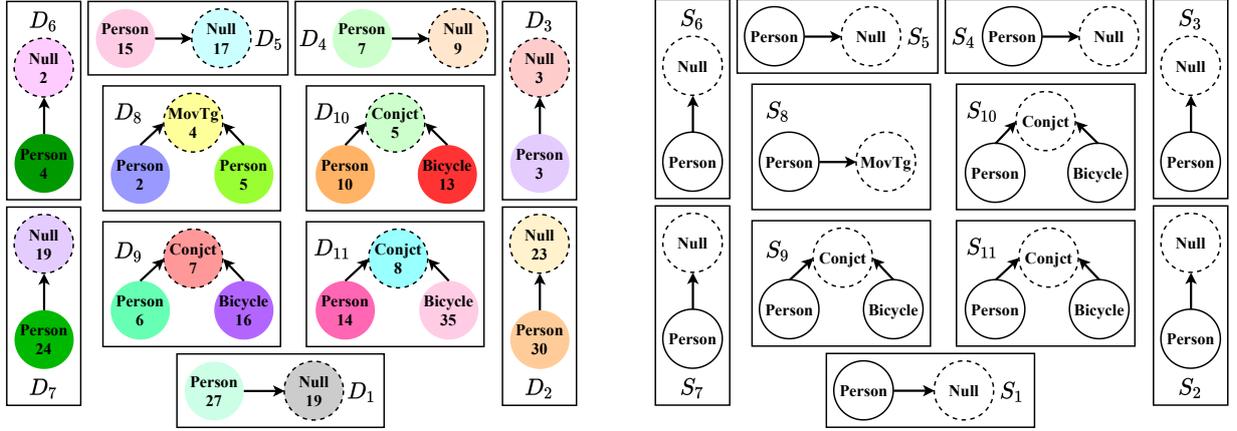Figure 18: Illustration of the instance-level and class-level graph representations.

22

(a) Instance-level graph splitting at post-processing. Graph $\mathcal{D}$ in Fig. 18b is split into its atomic components $D_m$ such that $\mathcal{D} = \{D_1, \ldots, D_{11}\}$. Corresponding attributes are split to atoms $\mathcal{A} = \{A_1, \ldots, A_{11}\}$ as well.

(b) Class-level abstractions are performed on disconnected subgraphs similar to the one in Fig. 18c so that we obtain $\mathcal{S} = \{S_1, \ldots, S_{11}\}$. Together with original graph $\mathcal{D}$ and attribute set $\mathcal{A}$, abstraction $\mathcal{S}$ forms the complete semantic description $Y$. Global goal-based filtering is performed using these high-level graphs.

Figure 19: Decomposing the instance-level and class-level graphs into their atomic components.

It should be noted that, while in Section 4, an instance level graph $\mathcal{D}$ is defined as union of disjoint subgraphs $\mathcal{D} = \{D_1, \ldots, D_m, \ldots D_M\}$, we use a single combined graph $D$ to represent the whole scene. It is desirable to split $\mathcal{D}$ into *atomic graphs* for an easy processing of goals, as discussed in Section 4. That is, splitting of $\mathcal{D}$ into its disjoint subgraphs $D_m$ (also called atomic graphs), can be performed during post-processing for this specific scenario. A similar splitting operation can also be applied to the attribute set $\mathcal{A} = \{A_1, \ldots, A_m, \ldots, A_M\}$ and class-level graph $\mathcal{S} = \{S_1, \ldots, S_m, \ldots, S_M\}$ where $S_m$ denotes the abstraction of $D_m$. The splitting operation on $\mathcal{D}$ and its class level counterpart is illustrated in Fig. 19.

We can also illustrate the goal-based filtering defined in Section 4 using the same scenario. Goal-based filtering operation allows for further distillation of the semantic information with respect to a specified goal and it enables reasoning over the semantic graph. In a sense, the filtering operation acts as a semantic parsing operator. The goal filtering operation for a small semantic graph with only two subgraphs is illustrated in Fig. 20. It should be observed that the queried graph patterns $G^S$ and $G^D$ as defined in (17), both include simple motifs with a single graph. Consequently, graph complexity and attribute complexity vectors reduce to scalars.

As illustrated in this part, the proposed goal-oriented semantic graph language and the signal processing framework can be incorporated into video signals for real-time applications. The semantic language allows for a structured and complete representation of the *meaningful* and *interesting* information embedded within the signal and allows for easy parsing of the information according to the desired goal.

## 5.2. Energy Efficient Sampling and Semantic Modeling of a Scalar Sensor Output

Energy efficiency is a critical research area for sensors, especially for the next generation of communication networks where the number of IoT devices is expected to rapidly increase. The proposed goal-oriented semantic language and the signal processing framework can enable dramatic reductions in transmitted information, especially when the sampled phenomenon is changing slowly with respect to the goals defined either internally or externally, leading to significant energy savings.

Another avenue of improvement for the efficiency of sensors can be introduced by employing smart non-uniform sampling strategies. For scalar sensors, if the signal output can be modeled as a random process whose statistical characterization is available, the sampling times can be adjusted according to desired *detection* and *missed detection* probabilities of certain events, e.g., a threshold crossing.
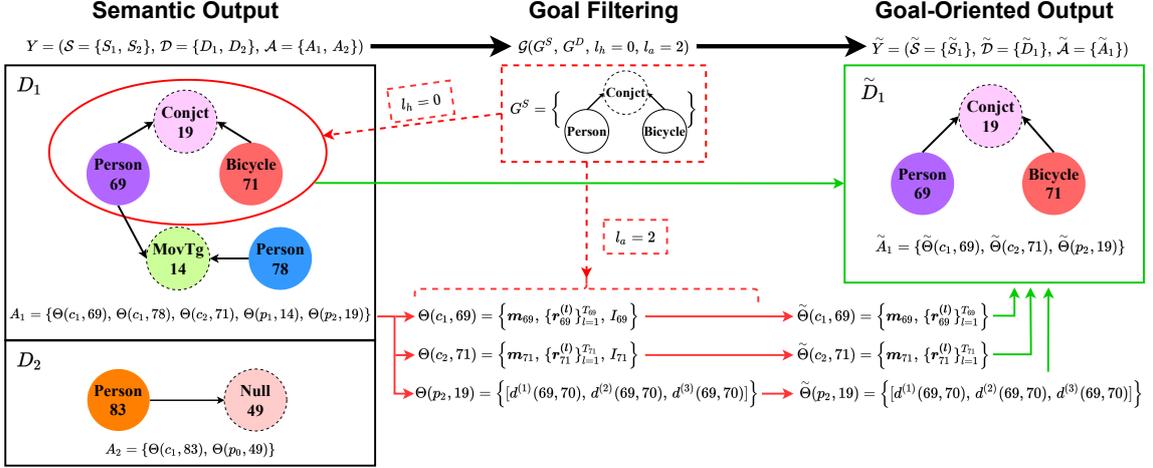
Figure 20: A goal-based filtering example. The filtering operation yields goal-oriented semantic description $\widetilde{Y}$. Dashed red lines denote the *queries/goals* and solid green lines denote the goal-oriented returns. For brevity, we do not visualize class-level graphs $S$. The goal pattern $G^S$ is searched in $S_1$ and $D_1$ to provide $\widetilde{S}_1$, and $\widetilde{D}_1$. Graph complexity is chosen as $l_h = 0$ so that only the exact matchings are returned. The corresponding attributes $A_1$ are distilled into $\widetilde{A}_1$ according to the given attribute complexity parameter $l_a = 2$.

In this subsection, we present a novel sampling strategy for a generic scalar sensor output, in conjunction with a goal-oriented semantic representation to improve the energy efficiency of sampling and transmission operations. First, we introduce a novel and optimal sampling strategy for scalar sensor outputs that conform to an auto-regressive AR($p$) Gaussian model, and present statistical results on the improved performance. Next, we implement the proposed goal-oriented semantic language definition in Section 4, to illustrate its validity even for relatively simple signal modalities.

*5.2.1. Optimal Sampling Strategies for Gaussian Auto-Regressive Models*

We start with a discrete time input signal $x_n$ over a horizon $H$, i.e., $n = 0 : H-1$, where $x_n$ is a Gaussian AR($p$) process with standard normal innovations and initial condition. For $p = 1$, we have an AR(1) process described by

$$x_n = \alpha x_{n-1} + \sqrt{1-\alpha^2} w_{n-1} \text{ for } n \geq 1, \tag{41}$$

where $x_0 \sim \mathcal{N}(0,1)$ and $w_i \sim \mathcal{N}(0,1)$. For a two-threshold event detection problem, we define $\tau \triangleq \inf \left\{ n > 0 : x_n \notin (l, u) \right\}$ where $u$, $l$ are the upper and lower thresholds, respectively. Assume that we take the first sample at time $n_1 = 1 + k$ for some non-negative $k$, and denote its probability of missing a threshold crossing as $P_{err}(k, l, u)$. We can relate this error probability with the CDF of $\tau$ as

$$
\begin{aligned}
P_{err}(k, l, u) &= \mathbb{P}(\tau \leq k) \\
&= 1 - \mathbb{P}\left( \left\{ \sup_{s \leq k} x_s < u \right\} \bigcap \left\{ \inf_{s \leq k} x_s > l \right\} \right) \\
&= 1 - \mathbb{P}\left( x_s \in (l, u) \quad \forall s \in \{1, 2, \ldots, k\} \right) \\
&= 1 - \psi_k, 
\end{aligned} \tag{42}
$$

where $\psi_k \triangleq \mathbb{P}\left( x_s \in (l, u) \quad \forall s \in \{1, 2, \ldots, k\} \right)$ with $\psi_0 = 1$ as a convention.
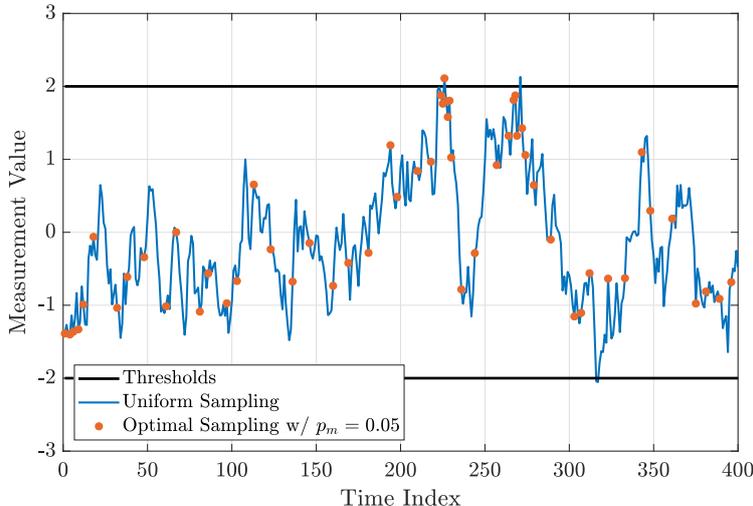
Figure 21: Optimal sampling strategy applied to an AR(1) process with two-thresholds.

Given a tolerable missed detection probability $p_m$, we choose the largest $k \in \mathbb{N}$ such that $\psi_k \geq 1 - p_m$ and take the first sample at $n_1 = 1 + k$. Since $x_n$ follows an AR($p$) process, calculation of $\psi_k$ can be done by using the Markov Property [165], i.e., factorization of the joint probability using conditional independence.

Next, we illustrate the derivation of the sampling algorithm for an AR(1) model. Generalization to a $p$-th order AR process can be performed similarly. We consider the model in (41) assuming that $x_0 \in (l, u)$ is observed. Therefore, we have $\mathbb{E}[x_n \mid x_0] = \alpha^n x_0$, $\mathbb{V}[x_n \mid x_0] = 1 - \alpha^{2n}$, and $\mathbb{C}[x_n, x_{n-1} \mid x_0] = \alpha(1 - \alpha^{2(n-1)})$, i.e., the mean is strictly decreasing towards 0 while the variance and covariance are strictly increasing towards 1 and $\alpha$, respectively.

To be able to calculate the next optimal sampling instance, we define the probability of the $i$-th sample being inside the two thresholds given that the previous sample was also inside the thresholds as

$$h_i \triangleq \mathbb{P}(x_i \in (l, u) \mid x_{i-1} \in (l, u)), \tag{43}$$

with an initial condition of $h_1 = \mathbb{P}(x_1 \in (l, u))$. Calculation of (43) requires only the first and second order statistics of $x_{i-1}$ and $x_i$, which can be computed numerically. Finally, the next optimal sampling time can be calculated as

$$n_1 = \inf\left\{n \in \mathbb{Z}_+ \,\middle|\, \psi_n < 1 - p_m\right\}, \tag{44}$$

where $\psi_k$ is defined using the Markov Property as

$$\psi_k = \prod_{i=1}^{k} h_i. \tag{45}$$

Once we observe the process at time $n_1$, we proceed to calculate the next sampling time $n_2$, and for any $k \in \mathbb{Z}^+$, the algorithm uses $x_{n_k}$ to determine the next optimal sampling time $n_{k+1}$.

In Fig. 21, we illustrate the proposed sampling algorithm for a first-order Gaussian AR process. Our aim is to detect crossings of thresholds $u = 2$ and $l = -2$. The desired maximum probability of missing the threshold crossings is set at 0.05. As expected, when the minimum distance to thresholds ($\min\{|x_n - l|, |u - x_n|\}$) decreases, the sampling rate increases.

We now investigate various statistical properties of (44) under the assumption

$$x_0 \sim \mathcal{N}(0, 1) \text{ is within the boundaries, i.e., } x_0 \in (l, u). \tag{46}$$
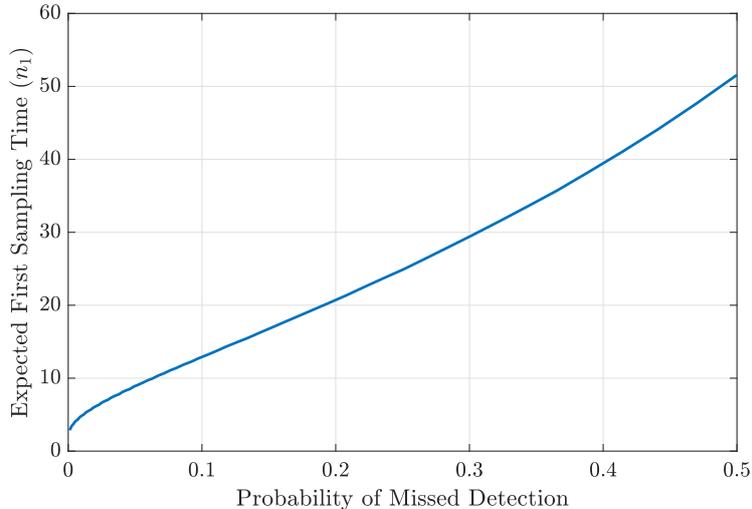
25

Figure 22: Expected first sampling time vs. probability of missed detection.

The distribution of the optimal sampling times requires calculation of (45), i.e., the products of dependent random variables $h_i$, under the assumption given in (46). The expectation of $n_1$ has significance since it represents the *average sampling period* of the algorithm. The expressions for the CDF and the expectation of $n_1$ are as follows:

$$\mathbb{P}\left(n_1 \leq k \,|\, x_0 \in (l, u)\right) = \mathbb{P}\left(\prod_{i=1}^{k} h_i < 1 - p_m\right), \tag{47}$$

$$\mathbb{E}\left[n_1 \,|\, x_0 \in (l, u)\right] = 1 + \sum_{k=1}^{\infty} \mathbb{P}\left(\prod_{i=1}^{k} h_i \geq 1 - p_m\right), \tag{48}$$

which can be computed numerically for given threshold levels $l, u$ and missed detection probability $p_m$.

Next, we present two simulation results. In the first simulation, we investigate the expected sampling period for various missed detection probabilities for AR(1) models. We use a first-order Gaussian AR process as in (41) with $\alpha = 0.95$, i.e.,

$$x_n = 0.95\, x_{n-1} + \sqrt{1 - 0.95^2}\, w_{n-1} \text{ for } n \geq 1, \tag{49}$$
$$x_0 \sim \mathcal{N}(0, 1) \text{ and } w_i \sim \mathcal{N}(0, 1).$$

Setting the thresholds $l = -2$, $u = 2$, we plot the expected first sampling time $\mathbb{E}\left[n_1 \,|\, x_0 \in (l, u)\right]$ versus the tolerated maximum probability of missed detections under the assumption that the current sample is within the boundaries in Fig. 22. As expected, the sampling period is positively correlated with the missed detection probability. The relationship is observed to be almost linear for this example.

In Fig. 23, we present the effect of the current sample value on the next sampling time interval for different missed detection probabilities $p_m$. The underlying assumptions and the model are the same as those in the previous simulation setup. We observe that as the minimum distance to the thresholds $(\min\{|x_n - l|, |u - x_n|\})$ increases, the constraint on the next sampling time that satisfies the missed detection probability requirement relaxes.

*5.2.2. Semantic Modeling Example for One-Dimensional Signals*

We can model the threshold detection problem studied in the previous subsection using the proposed semantic language introduced in Section 4. In a typical threshold detection problem, the semantic information extraction may include the detection of the region in which the signal resides at a given moment,
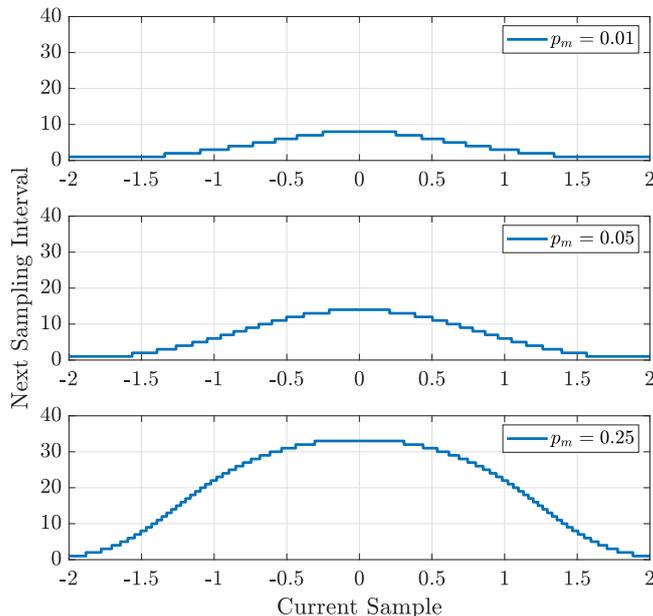
26

Figure 23: Next sampling time interval vs. current sample value for $p_m \in \{0.01, 0.05, 0.25\}$.

as well as the descriptions of threshold crossing events. Using the proposed formulation in Section 4, the components can be defined as the regions separated by the upper and lower thresholds that we denote as $u$ and $l$.

$$c_1 = \{x \in \mathbb{R} \mid x > u\} \tag{50}$$
$$c_2 = \{x \in \mathbb{R} \mid l \leq x \leq u\} \tag{51}$$
$$c_3 = \{x \in \mathbb{R} \mid x < l\} \tag{52}$$

We can easily extend the regions to a $k-$threshold model for $k > 2$ as well. Let $t_1 < t_2 < \ldots < t_k$ denote the thresholds in an ascending order. Then, the components can be defined as

$$c_1 = \{x \in \mathbb{R} \mid x > t_k\}, \tag{53}$$
$$c_i = \{x \in \mathbb{R} \mid t_{k-i+1} \leq x \leq t_{k-i+2}\}, \quad \forall i \in \{2, \ldots, k\}, \tag{54}$$
$$c_{k+1} = \{x \in \mathbb{R} \mid x < t_1\}. \tag{55}$$

The predicate set for this problem can be defined as descriptors of threshold crossing events as

$$P = (U, D, I), \tag{56}$$

where $U, D, I$ denote *upward crossings*, *downward crossings*, and *idling* (staying in the same region), respectively. Note that for this particular case, the *class multi-graph description* $\mathcal{S}_t$ and the *instance multi-graph description* $\mathcal{D}_t$ are defined as equivalent, since a higher level abstraction of a scalar signal is not required.

For the multi-level attribute set, a suitable choice for the first level is the current time index and amplitude of the process, and for each level $l$ we can store the previous $(l-1)$-th time and amplitude information as

$$\mathcal{A}_t = \{\Theta_{t,i}\} \text{ where } \Theta_{t,i} = \{\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(r)}\}, \tag{57}$$
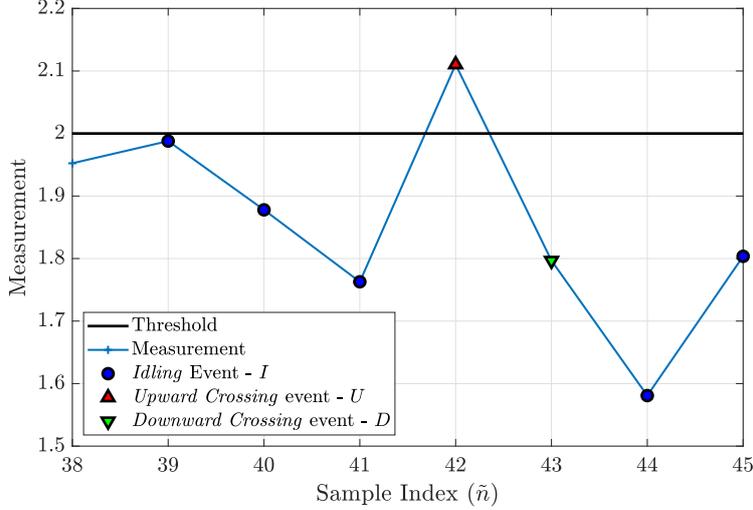$$\theta^l = (n - l + 1, x_{n-l+1}), \tag{58}$$

27

Figure 24: Illustration of *idling*, *upward crossing*, and *downward crossing* events for the AR(1) signal given in Fig. 21.

where $r$ is the total number of levels in the attribute sets. Depending on the specific application, alternative attribute levels including desired transform domain representations of the signal can be defined.

For a two-threshold case, and using the same Gaussian AR(1) process input signal in (49) with $l = -2$ and $u = 2$, we illustrate a region of interest around sampling indices $\tilde{n} \in [38, 45]$ in Fig. 24. Note that sampling points $\tilde{n}$ are not the same as the time index $n$, as we employ the optimal sampling strategy described in Section 5.2.1 to improve the energy and processing efficiency of the sensor. In Fig. 24, the samples marked with circles, upward-pointing triangles, and downward-pointing triangles represent *idling* ($I$), *upward crossing* ($U$), and *downward crossing* ($D$) events, respectively. In Fig. 25, class-level graph representations of the samples in Fig. 24, at $\tilde{n} = 41$ and $\tilde{n} = 42$, corresponding to $n = 225$ and $n = 226$ in Fig. 21, are illustrated with the previously defined semantic graph representation in (50)–(56). Additionally, for the semantic graphs shown in Fig. 25, the attribute sets as defined in (57),(58) for $r = 2$ are listed in Table 1.



(a) Semantic graph for $\tilde{n} = 41$ ($n = 225$).



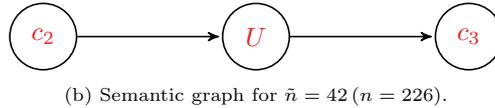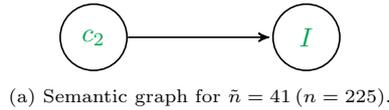(b) Semantic graph for $\tilde{n} = 42$ ($n = 226$).

Figure 25: Semantic graph descriptions for the events depicted in Fig. 24. (a) The *idling* event. (b) The *upward crossing* event.

Table 1: The attribute sets for samples taken for time instants 225 and 226, generated as companions to graphs given in Fig. 25.

|  | $\tilde{n} = 41$ ($n = 225$) | $\tilde{n} = 42$ ($n = 226$) |
|---|---|---|
| $\theta^1$ | (225, 1.76) | (226, 2.11) |
| $\theta^2$ | (224, 1.88) | (225, 1.76) |

As seen in Fig. 24, the signal is *idling* in Region-2 (between thresholds) at $\tilde{n} = 41$; hence, the corresponding component is $c_2$ and its predicate is $I$. At $\tilde{n} = 42$, the signal exceeds the upper threshold; therefore, the semantic description includes a flow from $c_2$ to $c_3$ with the predicate $U$ (*upward crossing*).
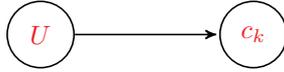
28

Figure 26: Goal pattern defined for an interest in detecting largest threshold crossings.

With the adoption of the semantic language in Section 4, we can easily define specific goals to reduce the storage or transmission rates for this particular application. In a $k-$threshold crossing detection problem, a typical goal might be the detection of exceeding the largest threshold, i.e., $x > t_k$. This means that we are interested in a *upward crossing* predicate connected to $c_k$, hence the goal pattern can be defined as shown in Fig. 26. With the goals as defined in Fig. 26, the transmission events (or, storage events for offline applications) can be reduced considerably. Coupled with the introduction of optimal sampling strategies in Section 5.2.1, the overall energy and processing efficiency of scalar sensors can be improved dramatically.

### 5.3. Discussion on the Applications of the Proposed Framework on Heterogeneous Multi-Sensor Networks

In this section, we presented detailed application examples using the proposed semantic signal processing framework on two extremes of signal complexity: namely, the real-time computer vision on video streams and 1-D scalar sensor outputs. The potentials of the proposed framework can be extended to many applications, including heterogeneous sensor networks that work on a dedicated task. A good example of the many fields that can benefit from semantic signal processing is agriculture.

For intelligent crop and plantation monitoring applications, a heterogeneous network of sensors (cameras, temperature/humidity sensors, etc.) provides information on critical events such as the crop yield and flowering status, temperature and humidity of the soil, existence of pests [166, 167]. In this application, a global goal regarding the ultimate objectives of the farm can be defined for a common customized language, and then it can be mapped to the respective capabilities of each sensor in the network. Based on goal-oriented returns from the sensors, appropriate actions can be taken. For example, an appropriate pesticide can be recommended and applied automatically to eliminate insects, or in the case of plant growth monitoring, ripe plants can be harvested or infected plants can be exterminated.

There are countless many other applications such as elderly fall detection, robot navigation, event detection in sports, animal monitoring, farm automation, traffic condition analysis, etc., that can employ the proposed framework or semantic approaches in general. We strongly believe that future research on signal processing should include a focus on different adaptations of a semantic framework for these applications.

## 6. Transmission of Goal-Oriented Semantic Information

The proposed semantic signal processing framework represents signals in a very organized and easy-to-parse structure, which enables goal-oriented filtering of the data to achieve very high compression rates. This is especially desirable in the next generation of machine-type communications where a huge amount of raw information will be generated by a plethora of IoT devices that needs to be transmitted throughout massive networks.

In this section, we first showcase the potential compression rates that can be achieved by semantic representation and goal-filtering through a simple example. Then, we discuss efficient coding and compression strategies for the proposed multi-graph and attribute set descriptions for storage and transmission, as well as strategies for transmission over noisy channels.

### 6.1. Data Compression Using Semantics and Goal-Filtering

Depending on the application and the nature of the signals of interest, the amount of data to be stored or transmitted can be greatly reduced using the proposed semantic signal processing framework. To illustrate this point via a simple example, we have simulated an object detection problem using YOLOv4 [18] on a 240-frame-long prerecorded video. The object set of YOLOv4 is defined as our component set $C$, while the predicate set is defined only with a $p_0$ : *exists* predicate, similar to the *null* predicate definition in Section 5.1. Therefore, in this example, each detection from a frame can be represented as a single component-predicate
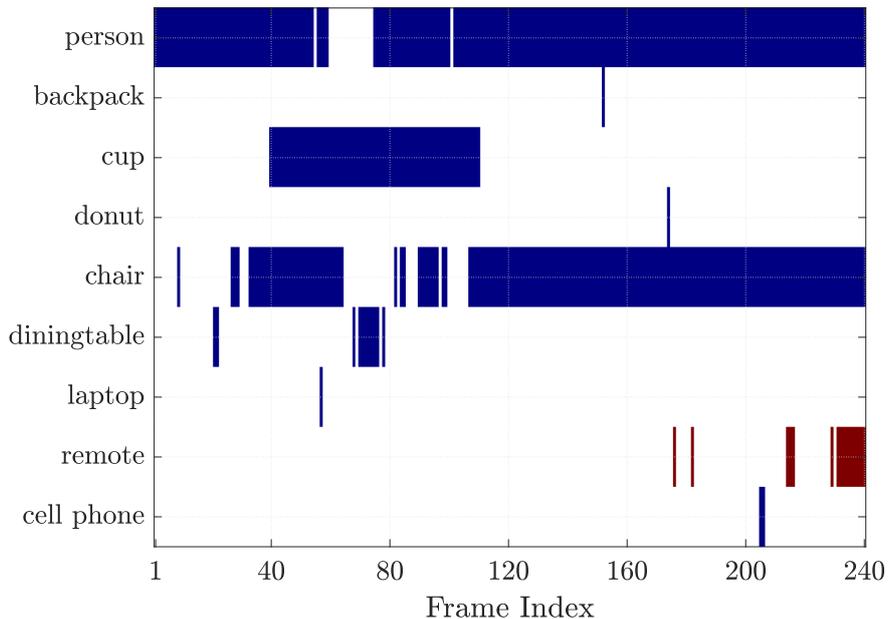
Figure 27: Detected component classes using YOLOv4 on a 240-frame long video. The *interesting* detections from class *remote* are shown in red.

connection (i.e., $c_i \rightarrow p_0$). The detected class instances throughout the 240-frame video are illustrated in Fig. 27.

In accordance with the semantic language definition in Section 4, the detector outputs in Fig. 27 correspond to a *multi-graph instance representation* $\mathcal{D}$. The corresponding attribute set $\mathcal{A}$ for this experiment is chosen as the cropped bounding box images of each detection, with a single level of complexity (i.e., $L = 1$). The class information is then encoded using Huffman coding [168] with a predefined historical occurrence rate. The bounding box images and the full frames are encoded using JPEG compression [169] with a constant compression rate of 10:1. Note that the specific coding schemes discussed here are selected only to give a general idea of the possible data throughputs, and alternative encoding schemes can be used for different applications.

A typical application for this object detection problem could be the transmission of *interesting* objects to an external agent. For the demonstration of the goal-filtering capabilities of the proposed framework, we assume that an external agent is interested in the detections of *remote* class, and may or may not require the bounding box images of pertinent detections.

With the above configuration of the experiment and the input video stream given in Fig. 27, the following transmission strategies are investigated:

- Full-Image Transmission: each full frame is sent,

- $(\mathcal{D}, \mathcal{A})$: the semantic graph outputs and the bounding box images are sent,

- $\mathcal{D}$: only the semantic graph outputs are sent,

- $(\widetilde{\mathcal{D}}, \widetilde{\mathcal{A}})$: the goal-filtered semantic graph outputs and the corresponding bounding box images are sent,

- $\widetilde{\mathcal{D}}$: only the goal-filtered semantic graph outputs are sent.

The data throughput per video frame using the above transmission strategies is given in Fig. 28. As seen in Fig. 28, the amount of data generated and transmitted can be reduced dramatically by organizing the data in a semantic framework (see $\mathcal{D}$, $\mathcal{A}$). Even further reductions are possible by introducing goals, and
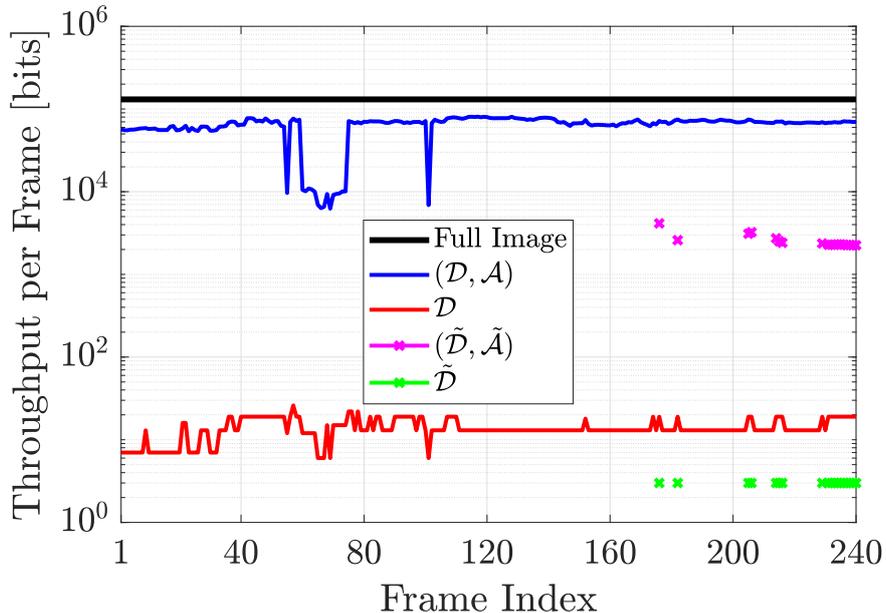
30

Figure 28: Data throughput using different transmission strategies.

filtering out unwanted signal components (see $\widetilde{\mathcal{D}}$, $\widetilde{\mathcal{A}}$). With a goal-oriented approach, transmissions are only sent when events of interest happen. Therefore, the amount of data reduction is dependent on the rate of innovation of *interesting* patterns in the signal. A summary of the results in Fig. 28 is given in Table 2, where the total throughput of the 240-frame video and the corresponding data compression rates are listed. As reinforced by Table 2, extremely high reductions in data rates are possible with the proposed goal-oriented approach.

Table 2: Total throughput and compression rates using the semantic representation of a video stream. Compression rates are given with respect to JPEG frame throughput.

|  | JPEG Frames | $(\mathcal{D}, \mathcal{A})$ | $\mathcal{D}$ | $(\widetilde{\mathcal{D}}, \widetilde{\mathcal{A}})$ | $\widetilde{\mathcal{D}}$ |
|---|---|---|---|---|---|
| **Total Throughput [bits]** | $3.1 \times 10^7$ | $1.6 \times 10^7$ | 3316 | $4.6 \times 10^4$ | 54 |
| **Compression Rate** |  | $2:1$ | $9468:1$ | $684:1$ | $580000:1$ |

As explained above, the encoding schemes used for this experiment are selected for demonstration purposes and ease of implementation. Throughout the rest of this section, we discuss efficient coding schemes to represent multi-graph representations with attribute sets, and possible transmission strategies for semantic information.

## 6.2. Efficient Transmission of Proposed Multi-Graph Representations and Attribute Sets

As described previously, graphs can be used to describe the goal-oriented semantic information acquired by different sensors. This description can be in the form of multi-graph class representations by using class atomic graphs $S_{t,i}$'s as given in (4), or more detailed information can be obtained through object multi-graph representations via object atomic graphs $D_{t,i}$'s as given in (11). As such, the storage and/or transmission of this graphical data warrant further discussion.

In order to store or transmit the graphical data, one may use the biadjacency matrix representation in (5). As extensively studied in the literature, an adjacency matrix can be described by $\mathcal{O}(n^2)$ bits where

$n$ is the number of vertices. The number of required bits can be reduced to $\mathcal{O}(n \log m + m \log n)$ by using adjacency list representation, where $m$ is the number of edges [170]. Since $n$ and $m$ can be reasonably small for practical scenarios, the resulting complexity becomes relatively low when the adjacency list representation is used. In addition, we can further compress the adjacency matrices by lossless compression techniques to increase the storage and transmission efficiencies. One way to compress graphs is to employ Huffman coding for the adjacency lists which is studied in [171] with a focus on web-graph structures. The compression rate of Huffman coding for this graph structure can be further increased by considering the similarities in shared links [172]. These studies can be easily extended to the proposed model in Section 4 to decrease the storage requirement and transmission time of the semantic description. However, Huffman coding uses the input probability distribution during the encoding process. Hence, we may need certain assumptions on the input probability distribution or implement a statistics gathering phase to efficiently compress the graphs in our language model.

Another line of research to compress the adjacency matrices includes universal codes for positive integers when the input probability distribution is not known as discussed in [170]. For instance, one may use Elias-$\gamma$ codes to encode the elements of the adjacency matrix which require $2\lceil \log x \rceil + 1$ bits to represent a positive integer $x$. This approach is especially preferred when the upper bound of the integers is not known beforehand [173]. This representation can be improved with Elias-$\delta$ coding with $\lceil \log x \rceil + 2\lceil \log \lceil \log x \rceil + 1 \rceil + 1$ bits, which is an asymptotically optimal universal code for positive integers [173].

The storage and transmission of adjacency matrices of graphical data require careful handling to utilize the communication resources effectively. Nonetheless, the communication and storage requirements in the proposed framework are dominated by the multi-level attribute sets of the nodes as defined in (14). A multi-level attribute set contains multiple levels of attributes. As previously discussed, an attribute set may carry information in varying levels of complexity, e.g., video stream applications may include raw image frames, extracted subfeatures of the components, along very simple scalar data. In the following, we focus on the transmission of high-rate data for semantic and goal-oriented communications.

Let us consider the $l$-th attribute of node-$j$ which can be a subfeature vector of the corresponding component in the scene graph as previously discussed in Section 5.1, and assume that $\boldsymbol{\theta}_{t,i}^{(l)}(n_j, k) \in \mathbb{R}^d$ for $l \in \{1, \cdots, L_{n_j}\}$ and $\boldsymbol{\theta}_{t,i}^{(l)}(n_j, k) \in [0, 1]$, where the dimension of the vector depends on the design of the semantic extractor. To store or transmit these subfeatures efficiently, we can apply lossy compression techniques to map each instance to a finite number of levels, i.e., we can perform quantization. One approach is to perform scalar quantization by addressing a single element of a vector at a time. Another approach is to employ vector quantization through a joint treatment of the entire vector at once, which is proven to be more effective than the scalar quantization [174]. A highly efficient and well-performing vector quantization algorithm is the Trellis Coded Quantization (TCQ) proposed in [175]. TCQ uses an expanded signal set and coded modulation with set partitioning. To transmit a signal from the codebook of size $2^m$, the original $2^m$-point constellations are doubled into $2^{m+1}$ points. A set partitioning applied to $2^{m+1}$ points to obtain $2^{\tilde{m}+1}$ subsets where $\tilde{m}$ is a positive integer which is less than $m$. Then, $\tilde{m}$ input bits are encoded by an $\tilde{m}/(\tilde{m}+1)$ rate convolutional code while the remaining bits are used to determine the codeword from the selected subset. Viterbi decoding [176] is employed to determine the sequence of codewords for a given input vector by minimizing the mean squared error (MSE) between the input and the output codeword. In TCQ, as depicted in Fig. 29, the paths through the trellis form the quantization codebook, and the path with the lowest cost is selected as the quantizer output for a given input sequence.

We aim at designing a lossy compression algorithm to store or transmit the semantic information with the smallest possible change on our inference. As a possible solution, cross-entropy can be considered as a cost metric since it quantifies the difference between two probability distributions. However, it is not easy to find the optimal quantizer output using the cross-entropy cost function for the given attribute, since the cost is not additive, which makes it unsuitable for use as a metric in Viterbi decoding. Instead, one may use other cost metrics such as *MSE, total variation distance (TVD), log-cosh cost,* and *quantile loss* defined as
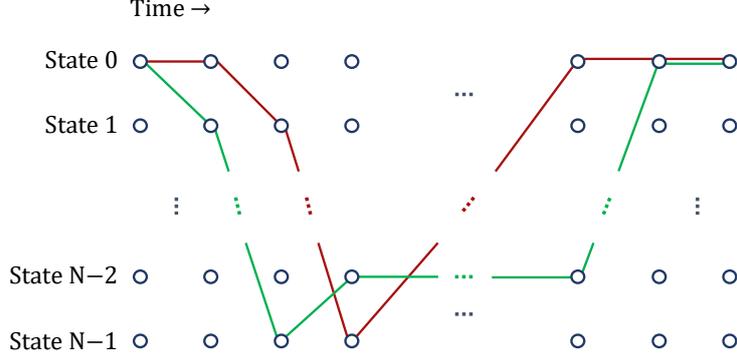
Figure 29: An illustrative trellis diagram. The path with the lowest cost is selected to determine the quantizer output.

- MSE cost:

$$L_{MSE} = \frac{1}{d} \sum_{b=1}^{d} \left( \theta_b - \hat{\theta}_b \right)^2,$$ (59)

- TVD cost:

$$L_{TVD} = \frac{1}{2} \sum_{b=1}^{d} \left| \theta_b - \hat{\theta}_b \right|,$$ (60)

- Log-cosh cost:

$$L_{LC} = \frac{1}{d} \sum_{b=1}^{d} \log \left( \cosh \left( \hat{\theta}_b - \theta_b \right) \right),$$ (61)

- Quantile loss with parameter $\gamma \in [0, 1]$:

$$L_Q = \frac{1}{d} \sum_{b=1}^{d} \max \left( \gamma \left( \hat{\theta}_b - \theta_b \right), (\gamma - 1) \left( \hat{\theta}_b - \theta_b \right) \right).$$ (62)

Note that for ease of notation, we drop the subscripts and superscripts denoting the parameters in the scene graph and denote the $b$-th element of the attribute vector and its quantized version by $\theta_b$ and $\hat{\theta}_b$, respectively, throughout this section. We also note that the TVD is a distance metric for probability distributions since the vectors to be compared are probability mass functions (PMFs) in the standard construction. In our study, we employ TVD for the subfeatures without any normalization, which proves to be a good metric for our purposes. Furthermore, as an extension to previously mentioned cost metrics, we also define the $l_p$-norm cost as

$$L_p = \frac{1}{d} \left( \sum_{b=1}^{d} \left( \theta_b - \hat{\theta}_b \right)^p \right)^{\frac{1}{p}},$$ (63)

which can be considered as a generalization of the MSE for $p = 1$, and the mean absolute error (MAE) for $p = 1$.

To summarize, while the construction of the proposed language to describe the semantic information is a significant foundation for goal-oriented signal processing, it is also crucial to address the efficient storage or transmission techniques of the generated semantic information. The adjacency matrices corresponding to the graph representations can be compressed by applying well-known methods in the literature, e.g., Huffman coding or universal codes. Furthermore, to decrease the bandwidth consumption of the dominant

components in the attribute sets, one may employ TCQ based quantization, where it is also possible to use different cost metrics to optimize the system performance.

### 6.2.1. A Simple TCQ Based Compression of Semantic Information

To demonstrate the performance of the semantic information compression, we use the same case study described in Section 5.1, where the sensors operate on video streams, and the language consists of 80 semantic components. To assess the performance of TCQ based quantization on a learning setup, we use the COCO dataset [160]. Employing the described CNN model in Section 5.1, the resulting feature vector of each detected object is stored as an element of the attribute set, i.e., $\boldsymbol{\theta} = \mathbf{r}_i \in \mathbb{R}^{128}$. These attributes can be stored locally or transmitted to a base station, where it is desired to protect semantic information (i.e., meaning) with a lossy compression scheme as much as possible. That is, we would like to compress the subfeature vector without changing the corresponding output of the classifier. To compress the subfeature vectors, we use TCQ based quantization with different cost metrics and a few bits, e.g., 2, 3, and 4 bits per element. Here, we do not follow Ungerboeck's construction for TCQ, which is described in [175]; instead, we employ an $m/(m+1)$ code rate convolutional code for $m$-bit quantization without set partitioning and map each vector entry to $m$ bits. For 2-bit quantization, a rate $2/3$ convolutional code with two connected shift register banks is used, where the constraint lengths of each shift register bank are three and two, respectively. In this architecture, the first output is connected to the first shift register array, the second one is connected to both shift register arrays, while the third output is connected to the last array only (see [177] for the shift register construction of convolutional codes). The corresponding octal code generator matrix specifying the output connections for each input bit is given by

$$G_{2/3} = \begin{bmatrix} 2 & 1 & 4 \\ 0 & 1 & 2 \end{bmatrix}. \tag{64}$$

Similarly, for the 3-bit vector quantization example, we use a rate $3/4$ convolutional code where the constraint lengths of register banks are $[5, 4, 4]$ with the generator matrix

$$G_{3/4} = \begin{bmatrix} 23 & 35 & 0 & 0 \\ 0 & 5 & 13 & 0 \\ 0 & 0 & 6 & 13 \end{bmatrix}. \tag{65}$$

Finally, for the 4-bit quantization example, rate $4/5$ convolutional code is employed where the constraint lengths of register banks are $[5, 4, 4, 4]$ with the generator matrix

$$G_{4/5} = \begin{bmatrix} 23 & 35 & 0 & 0 & 0 \\ 0 & 5 & 13 & 0 & 0 \\ 0 & 0 & 6 & 13 & 0 \\ 0 & 0 & 0 & 6 & 13 \end{bmatrix}. \tag{66}$$

Note that, in this simple case study, we do not optimize the architecture or the convolutional code used in TCQ. The optimization of this quantizer structure, particularly taking into account the nature of the semantic information can be a promising research direction.

The quantized feature vectors are fed into a neural network where we have two fully connected hidden layers with ReLU activations, and the softmax activation is used at the output layer. For comparison purposes, we also quantize the features by a uniform scalar quantizer which performs an element by element quantization. As a benchmark, we compare the TCQ-based quantizer's accuracy with a full resolution case in which the features are directly fed into the neural network. The accuracy of the benchmark is 83.12% where the elements of the feature vectors are represented by 32 bits. To evaluate the performance of the TCQ quantizer relative to the benchmark, we normalize the classification accuracies by the accuracy of the benchmark and report the results in Table 3. In this example, the best classification accuracy for 2-bit quantization is achieved with the Log-cosh cost while the MSE cost achieves higher accuracy than the other cost functions with 3-bit and 4-bit TCQ.

Table 3: Normalized classification accuracies of TCQ and uniform quantization where the accuracy of the benchmark is 83.12%, which is obtained by using a 32-bit representation of subfeatures. The normalization is performed by dividing the classification accuracy of the TCQ based semantic communication case study by the benchmark accuracy.

| Number of bits | Compression ratio | Uniform scalar quantization | TCQ | | | |
|---|---|---|---|---|---|---|
| | | | Cost metrics | | | |
| | | | MSE | Total variation distance | Log-cosh | Quantile $\gamma = 0.5$ |
| 2-bit | 0.06250 | 54.05% | 67.06% | 65.65% | 67.12% | 63.31% |
| 3-bit | 0.09375 | 86.15% | 91.87% | 91.21% | 91.86% | 90.80% |
| 4-bit | 0.12500 | 96.87% | 97.58% | 97.30% | 97.45% | 96.59% |

### 6.2.2. Further Extensions for Effective Compression

In many semantic signal processing applications, it is possible to further increase the compression rates. As an example, we consider the case study described in Section 5.1 where it is possible to track the individual objects by comparing their features as illustrated in Fig. 15 using the similarity of the feature vectors for consecutive frames. These similarities can also be exploited to create a differential data storage and transmit system. That is, the first feature vector obtained for a specific individual is to be compressed as previously described. When a new feature vector arrives, the difference with the previous one can be computed, which represents the fresh information in the new feature vector. One straightforward way to compress this differential information is to use TCQ with finer quantization levels suitable for the difference vectors, resulting in a higher resolution than sending all the feature vectors separately without utilizing the common information with the same strategy or communication costs. Furthermore, one may expect to have nearly *sparse* difference vectors due to similarities, and consequently, it may be possible to further compress the information to be stored or transmitted via innovative compressed sensing approaches [178].

The described goal-oriented semantic signal processing framework differs from traditional sensor networks in terms of the scope of knowledge generated and processed in sensors/base stations, which may alter the quantization/compression schemes. Unlike traditional sensor networks, the class information of each feature vector is already generated in the sensors. With this additional information, the quantizers can be designed for specific classes resulting in a *class-aware compression*. For instance, the transmit side can develop a performance indicator, e.g., using Kullback-Leibler distance and cross-entropy, to decide on the methods or metrics employed to compress the semantic information without losing the essential *meaning*. Such an indicator can be used to evaluate the performance of the TCQ-based quantizers with different cost metrics; thus, the sensor can decide on the quantizer output to be sent. Moreover, the averages of feature vectors for each class can be determined and shared with the receiver offline to enable *class-aware compression*. In most scenarios, the attribute set (14) contains the class information and its associated features. Hence, instead of compressing and transmitting each feature vector, sensors may use the average feature information which is available at both the sensors and the base station. Similar to a differential data transmission setup, the fresh information in the feature vector can be extracted by taking the difference between the new feature vector and the average feature vector of the corresponding class. Furthermore, vector quantization or compressed sensing can be employed to compress the differential feature vector. An interesting line of research in the context of *class-aware compression* is not only to use the class feature averages but also design a compressor for each class separately to exploit the shared knowledge of the sensors and the base station.

### 6.2.3. Transmission over Noisy Channels

As a further extension, joint source/channel coding resulting in more error-resilient systems can be employed [179] that are suitable for transmitting semantic information over noisy channels. In [180], the authors consider the joint design of channel-optimized vector quantization (COVQ) and rate-compatible punctured convolutional (RCPC) codes while [181] studies joint source and channel coding using multi-carrier modulation. The authors in [182] focus on image transmission where they analyze the rate-distortion behavior for joint source and channel coding with a wavelet-based sub-band source coding scheme and rate-compatible punctured convolutional (RCPC) code. Similar to these studies, a future line of research for

goal-oriented semantic communication may be to jointly optimize the source and channel coding where the desired meaning in the attribute set is transmitted to a base station in a robust manner despite the presence of channel noise.

As a concrete example of semantic communication over noisy channels, consider an $m$-bit quantization of the semantic information vector that is to be transmitted over a discrete memoryless channel (DMC). The $b$-th element of the subfeature vector $\boldsymbol{\theta} \in \mathbb{R}^d$ is quantized as $\hat{\boldsymbol{\theta}}$, where each element of $\hat{\boldsymbol{\theta}}$ can be represented by an $m$-bit binary sequence. As there are $2^m$ different $m$-bit sequences, let the sizes of both input and output alphabets of the DMC channel be $2^m$. Let us denote the $m$-bit input and output alphabets by $\mathcal{I}$ and $\mathcal{J}$, respectively, where $\mathcal{I} = \mathcal{J} = \{c_0, c_1, \cdots, c_{2^m-1}\}$. The channel transition probabilities are denoted by

$$\mathbb{P}(Y = c_i | X = c_j), \tag{67}$$

which corresponds to the probability of receiving the $m$-bit $c_j$ sequence when the $m$-bit $c_i$ sequence is transmitted. In this case, the relevant cost at the receiver side becomes

$$L_{DMC} = \sum_{b=1}^{d} \left( \sum_{k=0}^{2^m-1} \mathbb{P}(Y_b = c_k | X_b = x_b) l(x_b, c_k) \right), \tag{68}$$

where $\mathbb{P}(Y_b = c_k | X_b = x_b)$ is the transition probability for the $b$-th element of the subfeature vector representing the probability of receiving $m$-bit sequence $c_k \in \mathcal{J}$, when $x_b \in \mathcal{I}$ is transmitted, with the corresponding cost $l(x_b, c_k)$. Clearly, the overall cost function is assumed to be additive, and even though both $x_b$ and $c_k$ are $m$-bit binary sequences, the cost is calculated by considering their mapping into their real values, and the cost function can be one of the metrics described previously in this section. Note that $L_{DMC}$ is the average cost that takes into account the transition probabilities of the underlying channel model. As a further extension to the DMC channel model, one may also use different channel models, e.g., a Gaussian channel. In this case, instead of using the channel transition probabilities, the conditional probability density function defined by the channel should be used, and the summation over the variable $k$ in (68) should be replaced with an integral.

Note that while it is not explored here, joint optimization of the vector quantization with physical-layer properties to have an efficient and robust semantic communications system in practice may be an interesting research direction.

## 7. Conclusions and Future Research Directions

In this paper, for a structured and universal representation and efficient processing of the semantic information in signals, we propose and demonstrate a formal semantic signal processing framework in which the semantic information is represented by a graph-based structure. In the proposed framework, following a preprocessing stage of raw signals, a semantic information extractor identifies and classifies components from a set of predefined application-specific classes. The states, actions, and semantic relations among the identified components are described by another set of predefined application-specific predicates. Furthermore, along with the identification and classification of the components, each node in the graph is associated with a hierarchical set of attributes that provide additional information in an organized way. The proposed semantic signal processing framework enables the use of internally or externally defined *goals* that can vary with time. Using these goals, graphs can be grouped as those that are to be processed further and those that are of no interest. The further processing stages may include spatio-temporal tracking of the graph representations and a wide range of further signal processing operations on their more detailed set of attributes. At any point in the processing chain, the desired level of semantic information of those graphs which are of interest can be locally stored or shared with another processor through a communication protocol. Since a typical high bandwidth sensor data has sparse occurrences of interesting events, semantic signal processing of the sensor outputs results in remarkable compression rates. The wide range of applicability of the proposed goal-oriented semantic signal processing is illustrated over different examples and details are provided on how the proposed framework can be adapted for each use case.

The proposed semantic signal processing framework opens up multiple research directions in both theory and practice. First, available machine learning techniques should be assessed for their applicability in this framework for real-time, offline, and batch processing applications. For real-time applications where sensor data is semantically processed to observe and control the state of a system, semantic extraction should be completed within a certain deadline, depending on signal bandwidth and processing capabilities. Therefore, only those machine learning techniques that can extract semantic information with a low time complexity should be considered for this purpose. For applications that may allow offline processing such as semantic medical imaging, the semantic extraction and processing techniques must prioritize increasing the reliability of extracted meaning over the computational complexity. Based on the results of the assessment on both existing real-time and offline semantic extractors, desired features of new machine learning techniques should be identified for improved semantic extraction in the proposed framework.

Along with the new ML techniques for semantic information extraction, new goal-oriented signal processing techniques should be developed to take advantage of the available semantic information on the identified signal components. The proposed graph-based structure for the semantic information not only captures the relationship among different components through a set of application-specific predicates but also provides a hierarchical set of attributes for each component. Therefore, signal processing techniques that first identify the components of interest and then process their attributes can provide improved time complexity over conventional signal processing techniques, by prioritizing to achieve the desired accuracy of processing only on the components of interest. These algorithms should be able to set the goals of semantic filtering and semantic post-processing based on the estimated states of the components of interest. Further research on goal generation and dissemination of goals should also be performed. Mapping of the global goals defined by either human operators in a sensor network or defined by manufacturers of semantic devices at a hardware level, to local goals that can be interpreted by the pertinent device languages should also be studied.

Alternative hardware implementations of the proposed semantic signal processing framework should also be investigated. Since this framework can be implemented by using discrete subsystems or system-on-chip (SoC) architectures, a wide range of implementation alternatives exist. In the case of discrete subsystems, it is appropriate to combine the preprocessing and the semantic extraction stages in a sensor subsystem that generates data directly in the proposed graph-based structure with a hierarchical set of attributes. The aforementioned global goal definition by manufacturers can be employed in the semantic extraction stage at a hardware level, and the semantic filtering, semantic post-processing stages, and the storage unit can be built as a single subsystem. The transmission (or the storage) stage should be a versatile subsystem that can exchange relevant semantic data with a base station (or a storage unit). Detailed implementation of these subsystems should be carried out following an extensive analysis of their required specifications over multiple use cases. SoC implementation of the proposed framework will enable low-cost solutions for semantic signal processing and communication on a massive scale.

The introduction of the proposed semantic framework enables improvements over the limits imposed by the classical information theory, which is only concerned with the compression and transmission of bits as opposed to the meaning in the underlying data. In order to understand and quantify the new fundamental limits, modeling of the semantic noise, concept of information, relevant lossless compression, channel capacity and rate-distortion formulations have to be developed and computed in a variety of settings. In addition, as the goal-oriented approach also directly affects the aforementioned limits, the nature and effects of goals should be incorporated into these theoretical investigations as well. For efficient storage and communication of the desired semantic information with the proposed graph-based language model, it is essential to employ appropriate compression techniques which reduce the memory/bandwidth requirements. To compress the biadjacency matrix representations of the graph-based structure, well-known methods in the literature, e.g., Huffman coding or universal codes, can be employed. For the compression of the dominant components in the corresponding attribute sets, the use of TCQ-based quantization for these bandwidth-consuming vectors with different cost metrics is investigated in this paper; and these studies should be further extended. After identifying useful performance metrics, joint optimization of vector quantization with the physical-layer properties of the channel should be investigated for an efficient and robust system for semantic communications.

Finally, over the challenging use cases for the next generation of intelligent distributed sensor networks

and communication systems, the performance of the proposed semantic signal processing and communication methods should be compared with the available alternatives. The future of signal processing and communications requires a paradigm shift, and we believe that the proposed semantic framework may be a strong contender in developing the next generation of devices and systems whose performance surpass the limitations imposed by the classical information theory.

**Acknowledgements**

# References

[1] M. Kaveh, Signal processing everywhere [president's message], IEEE Signal Processing Magazine 28 (2) (2011) 6.

[2] Y. Bar-Hillel, R. Carnap, Semantic information, The British Journal for the Philosophy of Science 4 (14) (1953) 147–157.

[3] C. E. Shannon, A mathematical theory of communication, The Bell system technical journal 27 (3) (1948) 379–423.

[4] O. Javed, Z. Rasheed, K. Shafique, M. Shah, Tracking across multiple cameras with disjoint views, in: Proceedings of the Ninth IEEE International Conference on Computer Vision-Volume 2, 2003, p. 952.

[5] S. Khan, M. Shah, Consistent labeling of tracked objects in multiple cameras with overlapping fields of view, IEEE Transactions on Pattern Analysis and Machine Intelligence 25 (10) (2003) 1355–1360.

[6] K. S. Kumar, S. Prasad, P. K. Saroj, R. C. Tripathi, Multiple cameras using real time object tracking for surveillance and security system, in: 2010 3rd International Conference on Emerging Trends in Engineering and Technology, IEEE, 2010, pp. 213–218.

[7] E. I. Zacharaki, S. Wang, S. Chawla, D. Soo Yoo, R. Wolf, E. R. Melhem, C. Davatzikos, Classification of brain tumor type and grade using MRI texture and shape in a machine learning scheme, Magnetic Resonance in Medicine 62 (6) (2009) 1609–1618.

[8] J. Wiens, E. S. Shenoy, Machine learning for healthcare: On the verge of a major shift in healthcare epidemiology, Clinical Infectious Diseases 66 (1) (2017) 149–153.

[9] B. J. Erickson, P. Korfiatis, Z. Akkus, T. L. Kline, Machine learning for medical imaging, RadioGraphics 37 (2) (2017) 505–515.

[10] B. Krollner, B. Vanstone, G. Finnie, Financial time series forecasting with machine learning techniques: A survey, in: Proceedings of the 18th European Symposium on Artificial Neural Networks (ESANN 2010), 2010, pp. 25–30.

[11] A. Dingli, K. S. Fournier, Financial time series forecasting – a deep learning approach, International Journal of Machine Learning and Computing 7 (2017) 118–122.

[12] M. F. Dixon, I. Halperin, P. Bilokon, Machine Learning in Finance: From Theory to Practice, 1st Edition, Springer, Cham, 2020.

[13] Y. Zhou, D. Wilkinson, R. Schreiber, R. Pan, Large-scale parallel collaborative filtering for the Netflix prize, in: Algorithmic Aspects in Information and Management, Springer Berlin Heidelberg, Berlin, Heidelberg, 2008, pp. 337–348.

[14] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, et al., Grandmaster level in StarCraft II using multi-agent reinforcement learning, Nature 575 (7782) (2019) 350–354.

[15] N. Farsad, M. Rao, A. Goldsmith, Deep learning for joint source-channel coding of text, in: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2018, pp. 2326–2330.

[16] X. Wang, Y. Han, C. Wang, Q. Zhao, X. Chen, M. Chen, In-edge AI: intelligentizing mobile edge computing, caching and communication by federated learning, IEEE Network 33 (5) (2019) 156–165.

[17] J. Redmon, S. Divvala, R. Girshick, 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] C.-Y. Wang, A. Bochkovskiy, H.-Y. M. Liao, Scaled-yolov4: Scaling cross stage partial network (2021). `arXiv:2011.08036`.

[19] N. Tijtgat, W. Van Ranst, T. Goedeme, B. Volckaert, F. De Turck, Embedded real-time object detection for a UAV warning system, in: Proceedings of the IEEE International Conference on Computer Vision Workshops, 2017, pp. 2110–2118.

[20] L. Liu, H. Li, M. Gruteser, Edge assisted real-time object detection for mobile augmented reality, in: The 25th Annual International Conference on Mobile Computing and Networking, 2019, pp. 1–16.

[21] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, Y.-J. A. Zhang, The roadmap to 6G: AI empowered wireless networks, IEEE Communications Magazine 57 (8) (2019) 84–90.

[22] E. C. Strinati, S. Barbarossa, 6G networks: Beyond Shannon towards semantic and goal-oriented communications, Computer Networks 190 (2021) 107930.

[23] I. F. Akyildiz, A. Kak, S. Nie, 6G and beyond: The future of wireless communications systems, IEEE Access 8 (2020) 133995–134030.

[24] S. Dang, O. Amin, B. Shihada, M.-S. Alouini, What should 6G be?, Nature Electronics 3 (1) (2020) 20–29.

[25] Y. Chen, P. Zhu, G. He, X. Yan, H. Baligh, J. Wu, From connected people, connected things, to connected intelligence, in: 2020 2nd 6G wireless summit (6G SUMMIT), IEEE, 2020, pp. 1–7.

[26] C. E. Shannon, W. Weaver, The mathematical theory of information, Urbana: University of Illinois Press, USA, 1949.

[27] M. Schuster, K. K. Paliwal, Bidirectional recurrent neural networks, IEEE Transactions on Signal Processing 45 (11) (1997) 2673–2681.

[28] P. Simard, D. Steinkraus, J. Platt, Best practices for convolutional neural networks applied to visual document analysis, in: Proceedings. Seventh International Conference on Document Analysis and Recognition, Vol. 3, 2003, pp. 958–963.

[29] H. Larochelle, Y. Bengio, J. Louradour, P. Lamblin, Exploring strategies for training deep neural networks, Journal of Machine Learning Research 10 (1) (2009) 1–40.

[30] S. Hochreiter, J. Schmidhuber, Long short-term memory, Neural computation 9 (8) (1997) 1735–1780.

[31] Y. Li, W. Ouyang, B. Zhou, K. Wang, 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.

[32] D. Xu, Y. Zhu, C. B. Choy, 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.

[33] Y. Li, W. Ouyang, B. Zhou, J. Shi, C. Zhang, 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.

[34] J. Yang, J. Lu, S. Lee, D. Batra, D. Parikh, Graph R-CNN for scene graph generation, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 670–685.

[35] J. Gu, H. Zhao, Z. Lin, S. Li, J. Cai, M. Ling, Scene graph generation with external knowledge and image reconstruction, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 1969–1978.

[36] W. Wang, R. Wang, S. Shan, X. Chen, Exploring context and visual pattern of relationship for scene graph generation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 8188–8197.

[37] G. G. Chowdhury, Natural language processing, Annual review of information science and technology 37 (1) (2003) 51–89.

[38] D. L. Waltz, An English language question answering system for a large relational database, Communications of the ACM 21 (7) (1978) 526–539.

[39] L. Hirschman, R. Gaizauskas, Natural language question answering: the view from here, Natural Language Engineering 7 (4) (2001) 275.

[40] E. Choi, H. He, M. Iyyer, M. Yatskar, W.-t. Yih, Y. Choi, P. Liang, L. Zettlemoyer, Quac: Question answering in context, arXiv preprint arXiv:1808.07036.

[41] Q. You, H. Jin, Z. Wang, C. Fang, J. Luo, Image captioning with semantic attention, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 4651–4659.

[42] J. Aneja, A. Deshpande, A. G. Schwing, Convolutional image captioning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 5561–5570.

[43] Y. Pan, T. Yao, H. Li, T. Mei, Video captioning with transferred semantic attributes, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 6504–6512.

[44] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, J. Carlos Niebles, Dense-captioning events in videos, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 706–715.

[45] D. Marcu, The theory and practice of discourse parsing and summarization, MIT Press, 2000.

[46] R. Soricut, D. Marcu, Sentence level discourse parsing using syntactic and lexical information, in: Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics, 2003, pp. 228–235.

[47] J. Li, R. Li, E. Hovy, Recursive deep models for discourse parsing, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), 2014, pp. 2061–2069.

[48] M. Kountouris, N. Pappas, Semantics-empowered communication for networked intelligent systems, arXiv preprint arXiv:2007.11579.

[49] B. Juba, Universal Semantic Communication, Springer Berlin Heidelberg, 2011.

[50] J. Bao, P. Basu, M. Dean, C. Partridge, A. Swami, W. Leland, J. A. Hendler, Towards a theory of semantic communication, in: Proceedings of the 2011 IEEE 1st International Network Science Workshop, 2011, pp. 110–117.

[51] H. Xie, Z. Qin, A lite distributed semantic communication system for internet of things, IEEE Journal on Selected Areas in Communications 39 (1) (2020) 142–153.

[52] Z. Wang, J. Zhang, J. Feng, Z. Chen, Knowledge graph embedding by translating on hyperplanes, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 28, 2014, pp. 1112–1119.

[53] Q. Wang, Z. Mao, B. Wang, L. Guo, Knowledge graph embedding: A survey of approaches and applications, IEEE Transactions on Knowledge and Data Engineering 29 (12) (2017) 2724–2743.

[54] L. Zou, R. Huang, H. Wang, J. X. Yu, W. He, D. Zhao, Natural language question answering over RDF: a graph data driven approach, in: Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data, 2014, pp. 313–324.

[55] W. Zheng, H. Cheng, J. X. Yu, L. Zou, K. Zhao, Interactive natural language question answering over knowledge graphs, Information Sciences 481 (2019) 141–159.

[56] J. J. Carroll, I. Dickinson, C. Dollin, D. Reynolds, A. Seaborne, K. Wilkinson, Jena: Implementing the semantic web recommendations, in: Proceedings of the 13th International World Wide Web Conference on Alternate Track Papers & Posters, 2004, pp. 74–83.

[57] U. Bellur, R. Kulkarni, Improved matchmaking algorithm for semantic web services based on bipartite graph matching, in: IEEE International Conference on Web Services (ICWS 2007), IEEE, 2007, pp. 86–93.

[58] T. Segaran, C. Evans, J. Taylor, Programming the Semantic Web: Build Flexible Applications with Graph Data, ” O'Reilly Media, Inc.”, 2009.

[59] W. Fan, J. Li, S. Ma, N. Tang, Y. Wu, Y. Wu, Graph pattern matching: From intractable to polynomial time, Proc. VLDB Endow. 3 (1–2) (2010) 264—-275.

[60] W. Fan, X. Wang, Y. Wu, Incremental graph pattern matching, ACM Trans. Database Syst. 38 (3).

[61] F. Serratosa, Fast computation of bipartite graph matching, Pattern Recognition Letters 45 (2014) 244–250.

[62] J. Sun, H. Qu, D. Chakrabarti, C. Faloutsos, Neighborhood formation and anomaly detection in bipartite graphs, in: Fifth IEEE International Conference on Data Mining (ICDM'05), 2005, pp. 8 pp.–.

[63] Š. Čebirić, F. Goasdoué, H. Kondylakis, D. Kotzinos, I. Manolescu, G. Troullinou, M. Zneika, Summarizing semantic graphs: a survey, The VLDB Journal 28 (3) (2019) 295–327.

[64] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, Advances in Neural Information Processing Systems 25 (2012) 1097–1105.

[65] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 580–587.

[66] R. Girshick, Fast R-CNN, in: Proceedings of the IEEE Conference on Computer Vision, 2015, pp. 1440–1448.

[67] S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards real-time object detection with region proposal networks, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (6) (2016) 1137–1149.

[68] J. Redmon, A. Farhadi, YOLO9000: Better, faster, stronger, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 7263–7271.

[69] J. Redmon, A. Farhadi, YOLOv3: An incremental improvement, arXiv preprint arXiv:1804.02767.

[70] M. Tan, R. Pang, 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.

[71] W. Van Gansbeke, S. Vandenhende, S. Georgoulis, L. Van Gool, Unsupervised semantic segmentation by contrasting object mask proposals, arXiv preprint arXiv:2102.06191.

[72] X. Xia, B. Kulis, W-Net: A deep model for fully unsupervised image segmentation, arXiv preprint arXiv:1711.08506.

[73] X. Ji, J. F. Henriques, A. Vedaldi, Invariant information distillation for unsupervised image segmentation and clustering, arXiv preprint arXiv:1807.06653 2 (3) (2018) 8.

[74] J. Shotton, M. Johnson, R. Cipolla, Semantic texton forests for image categorization and segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2008, pp. 1–8.

[75] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, A. Blake, Real-time human pose recognition in parts from single depth images, in: CVPR 2011, IEEE, 2011, pp. 1297–1304.

[76] J. Tighe, M. Niethammer, 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.

[77] S. Hao, Y. Zhou, Y. Guo, A brief survey on semantic segmentation with deep learning, Neurocomputing 406 (2020) 302–321.

[78] F. Yu, V. Koltun, Multi-scale context aggregation by dilated convolutions, arXiv preprint arXiv:1511.07122.

[79] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3431–3440.

[80] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2015, pp. 234–241.

[81] H. Noh, S. Hong, B. Han, Learning deconvolution network for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision, 2015, pp. 1520–1528.

[82] V. Badrinarayanan, A. Kendall, R. Cipolla, Segnet: A deep convolutional encoder-decoder architecture for image segmentation, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (12) (2017) 2481–2495.

[83] R. P. Poudel, P. Lamata, G. Montana, Recurrent fully convolutional neural networks for multi-slice MRI cardiac segmentation, in: Reconstruction, Segmentation, and Analysis of Medical Images, Springer, 2016, pp. 83–94.

[84] P. Luc, C. Couprie, S. Chintala, J. Verbeek, Semantic segmentation using adversarial networks, arXiv preprint arXiv:1611.08408.

[85] Y. Xue, T. Xu, H. Zhang, L. R. Long, X. Huang, Segan: Adversarial network with multi-scale L-1 loss for medical image segmentation, Neuroinformatics 16 (3) (2018) 383–392.

[86] Y. Luo, Z. Zheng, L. Zheng, T. Guan, J. Yu, Y. Yang, Macro-micro adversarial network for human parsing, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 418–434.

[87] R. P. Poudel, S. Liwicki, R. Cipolla, Fast-SCNN: Fast semantic segmentation network, arXiv preprint arXiv:1902.04502.

[88] H. Park, L. L. Sjösund, Y. Yoo, J. Bang, N. Kwak, ExtremeC3Net: Extreme lightweight portrait segmentation networks using advanced C3-modules, arXiv preprint arXiv:1908.03093.

[89] H. Li, P. Xiong, H. Fan, J. Sun, Dfanet: Deep feature aggregation for real-time semantic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 9522–9531.

[90] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask R-CNN, in: Proceedings of the IEEE Conference on Computer Vision, 2017, pp. 2961–2969.

[91] A. Kirillov, K. He, R. Girshick, C. Rother, P. Dollár, Panoptic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 9404–9413.

[92] F. Lateef, Y. Ruichek, Survey on semantic segmentation using deep learning techniques, Neurocomputing 338 (2019) 321–348.

[93] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, P. Martinez-Gonzalez, J. Garcia-Rodriguez, A survey on deep learning techniques for image and video semantic segmentation, Applied Soft Computing 70 (2018) 41–65.

[94] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, 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.

[95] A. Karpathy, 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.

[96] J. Mao, W. Xu, Y. Yang, J. Wang, A. L. Yuille, Explain images with multimodal recurrent neural networks, arXiv preprint arXiv:1410.1090.

[97] X. Chen, 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.

[98] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, Y. Bengio, Show, attend and tell: Neural image caption generation with visual attention, in: International Conference on Machine Learning, PMLR, 2015, pp. 2048–2057.

[99] J. Johnson, A. Karpathy, 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.

[100] L. Yang, K. Tang, J. Yang, 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.

[101] D.-J. Kim, J. Choi, I. S. Kweon, et al., Dense relational image captioning via multi-task triple-stream networks, arXiv

preprint arXiv:2010.03855.

[102] 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 123 (1) (2017) 32–73.

[103] N. Aafaq, A. Mian, W. Liu, S. Z. Gilani, M. Shah, Video description: A survey of methods, datasets, and evaluation metrics, ACM Computing Surveys (CSUR) 52 (6) (2019) 1–37.

[104] S. Venugopalan, H. Xu, J. Donahue, M. Rohrbach, R. Mooney, K. Saenko, Translating videos to natural language using deep recurrent neural networks, arXiv preprint arXiv:1412.4729.

[105] H. Xu, S. Venugopalan, V. Ramanishka, M. Rohrbach, K. Saenko, A multi-scale multiple instance video description network, arXiv preprint arXiv:1505.05914.

[106] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, T. Darrell, K. Saenko, Sequence to sequence-video to text, in: Proceedings of the IEEE Conference on Computer Vision, 2015, pp. 4534–4542.

[107] H. Yu, J. Wang, Z. Huang, Y. Yang, W. Xu, Video paragraph captioning using hierarchical recurrent neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 4584–4593.

[108] J. Krause, J. Johnson, R. Krishna, 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.

[109] H. Xu, B. Li, V. Ramanishka, L. Sigal, 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.

[110] M. Z. Hossain, F. Sohel, M. F. Shiratuddin, H. Laga, A comprehensive survey of deep learning for image captioning, ACM Computing Surveys (CsUR) 51 (6) (2019) 1–36.

[111] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, L. Fei-Fei, Image retrieval using scene graphs, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3668–3678.

[112] X. Chang, P. Ren, P. Xu, Z. Li, X. Chen, A. Hauptmann, Scene graphs: A survey of generations and applications, arXiv preprint arXiv:2104.01111.

[113] Y. Li, W. Ouyang, X. Wang, 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.

[114] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, Y. Bengio, Graph attention networks, arXiv preprint arXiv:1710.10903.

[115] K. Tang, H. Zhang, B. Wu, W. Luo, 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.

[116] R. Zellers, M. Yatskar, S. Thomson, 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.

[117] R. J. Williams, Simple statistical gradient-following algorithms for connectionist reinforcement learning, Machine Learning 8 (3-4) (1992) 229–256.

[118] C. Lu, R. Krishna, M. Bernstein, L. Fei-Fei, Visual relationship detection with language priors, in: European Conference on Computer Vision, Springer, 2016, pp. 852–869.

[119] R. Yu, A. Li, V. I. Morariu, L. S. Davis, Visual relationship detection with internal and external linguistic knowledge distillation, in: Proceedings of the IEEE Conference on Computer Vision, 2017, pp. 1974–1982.

[120] A. Zareian, S. Karaman, S.-F. Chang, Bridging knowledge graphs to generate scene graphs, in: European Conference on Computer Vision, Springer, 2020, pp. 606–623.

[121] R. Wang, Z. Wei, P. Li, Q. Zhang, X. Huang, Storytelling from an image stream using scene graphs, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 9185–9192.

[122] H. Qi, Y. Xu, T. Yuan, T. Wu, S.-C. Zhu, Scene-centric joint parsing of cross-view videos, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 32, 2018, pp. 7292–7299.

[123] A. Agarwal, A. Mangal, et al., Visual relationship detection using scene graphs: A survey, arXiv preprint arXiv:2005.08045.

[124] M. Malik, M. K. Malik, K. Mehmood, I. Makhdoom, Automatic speech recognition: A survey, Multimedia Tools and Applications 80 (6) (2021) 9411–9457.

[125] K. H. Davis, R. Biddulph, S. Balashek, Automatic recognition of spoken digits, The Journal of the Acoustical Society of America 24 (6) (1952) 637–642.

[126] C.-C. Chiu, T. N. Sainath, Y. Wu, R. Prabhavalkar, P. Nguyen, Z. Chen, A. Kannan, R. J. Weiss, K. Rao, E. Gonina, et al., State-of-the-art speech recognition with sequence-to-sequence models, in: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2018, pp. 4774–4778.

[127] M. C. A. Korba, D. Messadeg, R. Djemili, H. Bourouba, Robust speech recognition using perceptual wavelet denoising and mel-frequency product spectrum cepstral coefficient features, Informatica 32 (3).

[128] R. Collobert, C. Puhrsch, G. Synnaeve, Wav2Letter: An end-to-end ConvNet-based speech recognition system, arXiv preprint arXiv:1609.03193.

[129] R. Polikar, et al., The wavelet tutorial (1996).

[130] M. Anusuya, S. Katti, Comparison of different speech feature extraction techniques with and without wavelet transform to kannada speech recognition, International Journal of Computer Applications 26 (4) (2011) 19–24.

[131] S. Ranjan, A discrete wavelet transform based approach to hindi speech recognition, in: International Conference on Signal Acquisition and Processing, IEEE, 2010, pp. 345–348.

[132] B. H. Juang, L. R. Rabiner, Hidden markov models for speech recognition, Technometrics 33 (3) (1991) 251–272.

[133] Ø. Birkenes, T. Matsui, K. Tanabe, S. M. Siniscalchi, T. A. Myrvoll, M. H. Johnsen, Penalized logistic regression with HMM log-likelihood regressors for speech recognition, IEEE Transactions on Audio, Speech, and Language Processing

18 (6) (2009) 1440–1454.

[134] H. Tang, C.-H. Meng, L.-S. Lee, An initial attempt for phoneme recognition using structured support vector machine (SVM), in: 2010 IEEE International Conference on Acoustics, Speech and Signal processing, IEEE, 2010, pp. 4926–4929.

[135] R. Solera-Ureña, J. Padrell-Sendra, D. Martín-Iglesias, A. Gallardo-Antolín, C. Peláez-Moreno, F. Díaz-de María, SVMs for automatic speech recognition: A survey, in: Progress in Nonlinear Speech Processing, Springer, 2007, pp. 190–216.

[136] S. E. Krüger, M. Schafföner, M. Katz, E. Andelic, A. Wendemuth, Speech recognition with support vector machines in a hybrid system, in: Ninth European Conference on Speech Communication and Technology, 2005, pp. 993–996.

[137] B. Wang, Y. Yin, H. Lin, Attention-based transducer for online speech recognition, arXiv preprint arXiv:2005.08497.

[138] J. Islam, M. Mubassira, M. R. Islam, A. K. Das, A speech recognition system for Bengali language using recurrent Neural network, in: 2019 IEEE 4th International Conference on Computer and Communication Systems (ICCCS), IEEE, 2019, pp. 73–76.

[139] A. Shewalkar, Performance evaluation of deep neural networks applied to speech recognition: RNN, LSTM and GRU, Journal of Artificial Intelligence and Soft Computing Research 9 (4) (2019) 235–245.

[140] T. Makino, H. Liao, Y. Assael, B. Shillingford, B. Garcia, O. Braga, O. Siohan, Recurrent neural network transducer for audio-visual speech recognition, in: IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), IEEE, 2019, pp. 905–912.

[141] D. Palaz, R. Collobert, et al., Analysis of CNN-based speech recognition system using raw speech as input, Tech. rep., Idiap (2015).

[142] S. Park, Y. Jeong, H. S. Kim, Multiresolution CNN for reverberant speech recognition, in: 2017 20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment (O-COCOSDA), IEEE, 2017, pp. 1–4.

[143] S. Ganapathy, V. Peddinti, 3-D CNN models for far-field multi-channel speech recognition, in: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2018, pp. 5499–5503.

[144] L. Besacier, E. Barnard, A. Karpov, T. Schultz, Automatic speech recognition for under-resourced languages: A survey, Speech Communication 56 (2014) 85–100.

[145] T.-Y. Tung, J. R. Pujol, S. Kobus, D. Gunduz, A joint learning and communication framework for multi-agent reinforcement learning over noisy channels, arXiv preprint arXiv:2101.10369.

[146] H. Xie, Z. Qin, G. Y. Li, B.-H. Juang, Deep learning enabled semantic communication systems, arXiv preprint arXiv:2006.10685.

[147] Z. Weng, Z. Qin, Semantic communication systems for speech transmission, arXiv preprint arXiv:2102.12605.

[148] J. Hu, L. Shen, G. Sun, Squeeze-and-excitation networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 7132–7141.

[149] B. Güler, A. Yener, A. Swami, The semantic communication game, IEEE Transactions on Cognitive Communications and Networking 4 (4) (2018) 787–802.

[150] G. Shi, D. Gao, X. Song, J. Chai, M. Yang, X. Xie, L. Li, X. Li, A new communication paradigm: from bit accuracy to semantic fidelity, arXiv preprint arXiv:2101.12649.

[151] 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.

[152] D. Gündüz, D. B. Kurka, M. Jankowski, M. M. Amiri, E. Ozfatura, S. Sreekumar, Communicate to learn at the edge, IEEE Communications Magazine 58 (12) (2020) 14–19.

[153] L. Cohen, Time-frequency analysis, PTR Prentice Hall, 2012.

[154] Y. Chan, Wavelet basics, Springer Science & Business Media, 1994.

[155] C. Torrence, G. P. Compo, A practical guide to wavelet analysis, Bulletin of the American Meteorological society 79 (1) (1998) 61–78.

[156] A. K. Özdemir, S. Karakaş, E. D. Çakmak, D. İlhan Tüfekçi, O. Arıkan, Time–frequency component analyser and its application to brain oscillatory activity, Journal of Neuroscience Methods 145 (1) (2005) 107–125.

[157] R. E. Hall, B. Bowerman, J. Braverman, J. Taylor, H. Todosow, U. Von Wimmersperg, The vision of a smart city, Tech. rep., Brookhaven National Lab., Upton, NY (US) (2000).

[158] K. Su, J. Li, H. Fu, Smart city and the applications, in: International Conference on Electronics, Communications and Control (ICECC), IEEE, 2011, pp. 1028–1031.

[159] N. Wojke, A. Bewley, 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.

[160] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, P. Dollár, Microsoft COCO: common objects in context (2015). `arXiv:1405.0312`.

[161] X. Du, T.-Y. Lin, P. Jin, G. Ghiasi, M. Tan, Y. Cui, Q. V. Le, X. Song, Spinenet: Learning scale-permuted backbone for recognition and localization, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 11592–11601.

[162] G. Welch, G. Bishop, et al., An introduction to the Kalman filter, Chapel Hill, NC, USA, 1995.

[163] N. Wojke, A. Bewley, Deep cosine metric learning for person re-identification, in: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV), IEEE, 2018, pp. 748–756.

[164] R. Jonker, T. Volgenant, Improving the hungarian assignment algorithm, Operations Research Letters 5 (4) (1986) 171–175.

[165] A. A. Markov, The theory of algorithms, Trudy Matematicheskogo Instituta Imeni VA Steklova 42 (1954) 3–375.

[166] A. Baggio, Wireless sensor networks in precision agriculture, in: ACM Workshop on Real-World Wireless Sensor Networks (REALWSN 2005), Stockholm, Sweden, Vol. 20, Citeseer, 2005, pp. 1567–1576.

[167] M. Srbinovska, C. Gavrovski, V. Dimcev, A. Krkoleva, V. Borozan, Environmental parameters monitoring in precision agriculture using wireless sensor networks, Journal of Cleaner Production 88 (2015) 297–307.

[168] D. A. Huffman, A method for the construction of minimum-redundancy codes, Proceedings of the IRE 40 (9) (1952) 1098–1101.

[169] W. B. Pennebaker, J. L. Mitchell, JPEG: Still image data compression standard, Springer Science & Business Media, 1992.

[170] M. Besta, T. Hoefler, Survey and taxonomy of lossless graph compression and space-efficient graph representations, arXiv preprint arXiv:1806.01799.

[171] T. Suel, J. Yuan, Compressing the graph structure of the web, in: Proceedings DCC 2001. Data Compression Conference, 2001, pp. 213–222.

[172] M. Adler, M. Mitzenmacher, Towards compressing web graphs, in: Proceedings DCC 2001. Data Compression Conference, 2001, pp. 203–212.

[173] P. Elias, Universal codeword sets and representations of the integers, IEEE Transactions on Information Theory 21 (2) (1975) 194–203.

[174] T. M. Cover, J. A. Thomas, Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing), Wiley-Interscience, USA, 2006.

[175] M. W. Marcellin, T. R. Fischer, Trellis coded quantization of memoryless and gauss-markov sources, IEEE Transactions on Communications 38 (1) (1990) 82–93.

[176] N. Seshadri, C.-E. Sundberg, List viterbi decoding algorithms with applications, IEEE Transactions on Communications 42 (234) (1994) 313–323.

[177] W. Ryan, S. Lin, Channel codes: Classical and modern, Cambridge University Press, 2009.

[178] R. G. Baraniuk, Compressive sensing [lecture notes], IEEE Signal Processing Magazine 24 (4) (2007) 118–121.

[179] C. E. Shannon, Communication in the presence of noise, Proceedings of the IRE 37 (1) (1949) 10–21.

[180] A. J. Goldsmith, M. Effros, Joint design of fixed-rate source codes and multiresolution channel codes, IEEE Transactions on Communications 46 (10) (1998) 1301–1312.

[181] K.-P. Ho, J. M. Kahn, Transmission of analog signals using multicarrier modulation: A combined source-channel coding approach, IEEE Transactions on Communications 44 (11) (1996) 1432–1443.

[182] M. J. Ruf, J. W. Modestino, Rate-distortion performance for joint source and channel coding of images, in: Proceedings., International Conference on Image Processing, Vol. 2, IEEE, 1995, pp. 77–80.