

# AI Oriented Large-Scale Video Management for Smart City: Technologies, Standards and Beyond

Lingyu Duan<sup>1</sup>, Yihang Lou<sup>1,2</sup>, Shiqi Wang<sup>3</sup>, Wen Gao<sup>1,2</sup>, *Fellow, IEEE*, Yong Rui<sup>4</sup>, *Fellow, IEEE*

<sup>1</sup>Institute of Digital Media, Peking University, Beijing, China

<sup>2</sup>SECE of Shenzhen Graduate School, Peking University, Shenzhen, China

<sup>3</sup>Department of Computer Science, City University of Hong Kong

<sup>4</sup>Lenovo Research, Beijing, China



**Abstract**—Deep learning has achieved substantial success in a series of tasks in computer vision. Intelligent video analysis, which can be broadly applied to video surveillance in various smart city applications, can also be driven by such powerful deep learning engines. To practically facilitate deep neural network models in the large-scale video analysis, there are still unprecedented challenges for the large-scale video data management. Deep feature coding, instead of video coding, provides a practical solution for handling the large-scale video surveillance data. To enable interoperability in the context of deep feature coding, standardization is urgent and important. However, due to the explosion of deep learning algorithms and the particularity of feature coding, there are numerous remaining problems in the standardization process. This paper envisions the future deep feature coding standard for the AI oriented large-scale video management, and discusses existing techniques, standards and possible solutions for these open problems.

**Index Terms**—Video analysis, deep learning, smart city, deep feature coding standard.

## 1 INTRODUCTION

Recently, a considerable number of deep learning algorithms have been proposed, which exhibit substantial performance improvement in various computer vision tasks. Compared with traditional handcrafted features, deep learning algorithms aim to learn representative features from the vast amounts of training data. Since AlexNet [1] won the ImageNet competition, there are tremendous research activities focusing on designing more powerful and deeper networks. Follow ups like VGGNet [2], GoogleNet [3], ResNet [4] and DenseNet [5] have greatly improved the discrimination capability of features to a higher level, which also boosted the performance of many visual analysis tasks. Generally speaking, these technologies have naturally made substantial impact on public security, such as face recognition [6], person [7] and vehicle reidentification [8] in surveillance videos.

Recent years have witnessed dramatically increased demand for the smart city construction, where the concerning safety issues have received sufficient interest. In particular,

there is a vast and increasing proliferation of surveillance videos acquired and transmitted over both wireline and wireless networks. Due to the real-time recording of the physical world, surveillance video is very valuable and there is considerable concern regarding how to efficiently manage such surveillance video big data. In view of the explosion of the surveillance systems deployed in urban areas and millions of objects/events captured every day, there are a unique set of challenges regarding efficient analysis and search. In particular, video compression and transmission constitute the basic infrastructure to support these applications. Though the state-of-the-art video coding standards such as H.265/HEVC have dramatically improved the coding performance, it is still questionable that whether such big video data can be efficiently handled by visual signal level compression. Fortunately, an alternative strategy “analyze then compress” provides a solution, which transmits the compact features extracted and compressed at the edge end to the server side. Such paradigm can sufficiently satisfy various intelligent video analysis tasks, by using significantly less data than the compressed video itself. In Fig. 1, the infrastructure of the smart city with large-scale video management based on feature extraction and transmission is illustrated. In particular, to meet the demand of large-scale video analysis in smart city applications, the feature stream instead of video signal stream can be transmitted. As such, the intelligent front end devices extract features locally and then convey the encoded feature stream to the server for analysis purpose.

While the field of artificial intelligence is still quickly evolving, and efficient and novel deep learning algorithms will continue to emerge in the coming years, it is also interesting to discuss how we could enable the interoperability of the compressed deep learning features in real-world applications. In contrast with video coding, which directly compresses the visual signals into the bitstream, feature coding involves both feature extraction and compression process. In particular, feature extraction serves as the raw features producer to generate the source for compression

and is responsible for the answer of what to compression. Feature compression accounts for the conversion of raw deep features into compact representation bitstream. The purpose of this article is to provide an overview of the existing deep learning techniques in video surveillance and envisioning the future deep learning feature coding standards. We will start by a brief review of the current status of deep learning in video surveillance, followed by discussions on the compact feature standard in MPEG. Then the open problems of deep feature coding standardization will be discussed, where we can perceive both great promises and challenges.

## 2 STANDARDIZATION OF COMPACT HAND-CRAFTED FEATURE DESCRIPTOR

### 2.1 Compact Descriptors for Visual Search

In view of the importance of the transmission of feature descriptors, MPEG has finalized the standardization of Compact Descriptors for Visual Search (CDVS) [9] and published the standard in Sep. 2015. In CDVS, handcrafted local and global descriptors are leveraged to represent the visual characteristics of images. The normative blocks of CDVS involve the extraction of local and global descriptors. More specifically, the local descriptors consist of SIFT descriptors which are efficiently compressed by a low-complexity transform coding. The raw local descriptors are further selected and aggregated to generate a Scalable Compressed Fisher Vector (SCFV), with competitive matching accuracy and low memory footprint. In view of the fluctuation of available bandwidth in the mobile environment, CDVS supports interoperability between different size bitstream by setting six operating points from 512B to 16KB. More technical details about CDVS are referred to [9].

### 2.2 Compact Descriptors for Video Analysis

The emerging requirements of video analysis facilitate the standardization of large-scale video analysis. The MPEG has moved forward to standardize Compact Descriptors for Video Analysis (CDVA) [10]. Video consists of sequentially correlated frames, such that extracting the feature from each frame leads to high redundancy in feature representation and unnecessary computational costs. The ongoing CDVA standard adopted multi-keyframe based image retrieval, which converts the problem of video retrieval into an image retrieval task. In particular, the local and global descriptors of standardized CDVS descriptors are extracted on the sampled keyframes of a given query video, which are further packed together to constitute the CDVA descriptors. Moreover, the deep learning features [11] are also adopted to further boost the analysis performance. Fig. 2 presents the framework of ongoing CDVA with handcrafted features and deep features, including the video structure and the normative components of feature extraction. It is also worth mentioning that the NIP descriptor has been adopted into the working draft of CDVA standard.

## 3 DEEP LEARNING IN VIDEO SURVEILLANCE

With the exponential growth of the video surveillance data, video content analysis has been a long standing research

Table 1  
The core techniques involved in the visual system of smart city.

The Core Techniques	Description
Feature Generation	Extract discriminative feature representation
Feature Generalization	Generalize the deep feature, to different tasks e.g., from person ReID to vehicle ReID.
Feature Redundancy Removal	Remove the redundancies of the features in spatial or temporal domain
Rate-Distortion Optimization	Investigate the distortion of features and optimize feature compression with rate-distortion optimization.
Feature Binarization	Binarize features to enable fast feature transmission and analysis.
Network Compression	Efficiently represent the network and lower the disk and memory cost.

topic in computer vision community. Currently, there are four urgent visual analysis tasks in the surveillance scenario, i.e., image/video retrieval, person Re-Identification, face recognition and vehicle retrieval. These tasks play important roles in building the safety city and ensuring the public security. Associated with these tasks, the core techniques in establishing the visual system of the smart city are shown in Table 1. Multiple academic disciplines, including visual signal processing, computer vision, compression as well as hardware architectures are involved in the further construction of the system. It is envisioned that with the standardization of deep learning features and the advancements of these technologies, the system that sees intelligently, efficiently and greenly will eventually come true.

### 3.1 Image/Video retrieval

Image/video retrieval refers to searching for the images/videos representing the same objects or scenes as the one depicted in the query, which may present under different scales, illuminations, rotations or even occlusions. In the last decade, the image/video retrieval has benefited a lot from handcrafted SIFT descriptors due to its robustness to the image transformations. However, after the AlexNet [1] won ILSVRC12 by a significant margin, the CNN-based image feature representation has become mainstream techniques when handling complex and semantic vision analysis. Regarding to both image and video retrieval, competitive and even better retrieval performance [12] [12] [13] [11] has been reported on several benchmarks.

CNN-based retrieval methods can be categorized into two types: pre-trained and fine-tuned CNN models. The commonly used pre-trained CNN models are trained on ImageNet dataset consisting of 1.2 million images of 1000 classes, such that the features can be regarded as generic. The descriptors can be extracted from fully-connected (FC) layers or intermediate layers. The FC descriptors have a global receptive field and the intermediate local descriptors have a smaller receptive field and location information encoded in 2D feature maps. To obtain the global representation, encodings like VLAD and FV are usually adopted. In addition, the direct pooling can also generate discriminative features. For example, in [13] Toliás et al. employed max pooling on selected regions in intermediate feature maps and subsequently performed sum pooling. Though impressive results can be achieved by pre-trained model, there is a

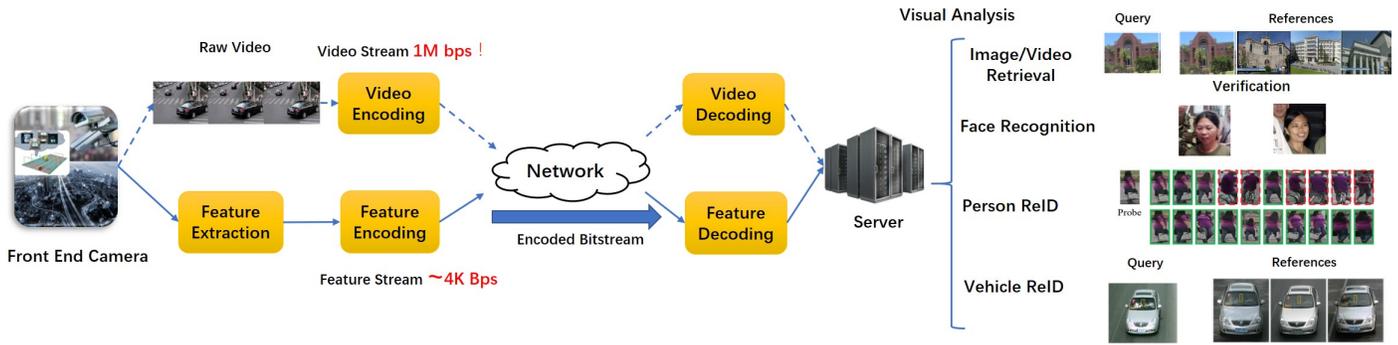


Figure 1. The infrastructure of the large-scale video management with feature transmission for smart city applications.

trend to fine-tune CNN model on a task-oriented dataset for specific retrieval. The classification and verification based networks are two typical types. The former is trained to classify pre-defined categories, and later adopts siamese network [14] with contrastive loss or triplet loss. On several retrieval benchmarks such as Holidays, Oxford5K, Paris 6K, the fine-tuned models have achieved the state-of-the-art performance.

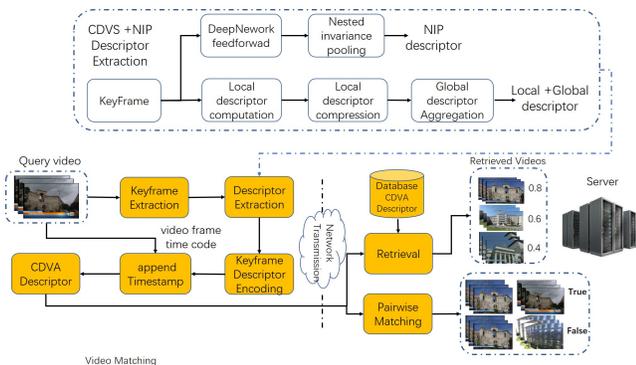


Figure 2. Illustration of MPEG CDVA evaluation framework

### 3.2 Person Reidentification

Person ReIdentification (ReID) has attracted more and more research focus due to its application significance in video surveillance. It aims to search whether a given person is present in other cameras. The widespread camera deployment in public places and the increased safety requirements make the existing manual labor spotting scheme powerless when facing the real-time generated massive video data. A practical person ReID system involves person detection, tracking, and retrieval. In particular, person retrieval is the main research focus among the related works. The challenge of this task is how to accurately match two images of the same person under variant scenes, viewpoints, scales and lighting conditions.

Deep ReID systems mostly employ two types of CNN models, i.e., siamese and classification models. The difference between these two models lies in the input form and definition of the loss function. The siamese models [15] [16] leverage image pairs or triplets as input, then let them forward propagate to get feature vector in embedding space. The distance of images with the same person is constrained

on the feature vector by a minimum margin using contrastive loss or triplet loss. By contrast, the classification models [17] [18] treat each person identity as a class, and the classification based loss functions are usually employed, such as softmax loss. These models pay more attention to the feature representation from the perspective of local and global combination, part-attention model, human body's skeleton model, etc. Intuitively, the combinations of siamese and classification model have also received a lot of attention.

### 3.3 Face recognition

Due to the nonintrusive recognition manner (intrusive like fingerprint, retina recognition), face recognition has great application potential in surveillance security. Over the last decades, amount of works in face recognition [6] [19] have emerged, which greatly boosted the accuracy on the popular benchmark such as Labeled Face in the Wild (LFW) to an unprecedented level. In real-world applications, the captured face images may not be as high quality as that in LFW dataset, creating many challenging problems originated from arbitrary poses, low quality resolutions, occlusions and small scales.

Typically, face recognition includes face identification and face verification. The former classifies a given face as a specific identity, and the latter verifies whether a given face pair belongs to the same identity. Regarding experimental setup, there are closed-set and open-set settings. Under closed-set, the testing identities are contained in the training set. By contrast, in open-set, the testing identities do not appear in the training set. Therefore, the real-world recognition can be regarded as face verification in open-set. Essentially, this task is defined as a metric learning problem. The expected feature representation should be able to meet the demand of small intra-class distance and large inter-class distance. Deep models are capable of building the above criterion by setting appropriate loss functions.

### 3.4 Vehicle Reidentification

In many vehicle-relevant tasks, vehicle ReIdentification is the most crucial technique in city security. The license plate is usually the straightforward choice to identify a vehicle. However, in real applications, most surveillance cameras are not equipped with recognition capability. Furthermore, the license plates of vehicles in many cases are occluded or faked. Thus, the visual appearance based techniques present

Table 2

Performance comparisons of methods reported on CDVA benchmarks where the landmarks, scene and objects are the sub datasets

Methods	Dims	Year	Land mark	Scene	Objects	All
CXM0.2 [21]	1024	2016	0.598	0.594	0.917	0.721
MAC [13]	512	2015	0.619	0.762	0.718	0.670
SPoC [12]	256	2015	0.691	0.840	0.703	0.709
CroW [12]	512	2015	0.639	0.784	0.720	0.683
R-MAC [13]	512	2015	0.746	0.873	0.782	0.771
HNIP VGG [11]	512	2016	0.748	0.901	0.850	0.801
HNIP Alex [11]	768	2016	-	-	-	0.772
HNIP Res [11]	2048	2016	-	-	-	0.817

Table 3

Performance comparisons of vehicle re-identification methods on VehicleID and VeRI benchmarks

Methods	Dims	Year	VEHICLEID	VERI-776
Triplet Loss [22]	400	2015	0.373	-
Softmax Loss	1024	2015	0.580	0.343
Triplet+Softmax Loss [20]	1024	2014	0.650	0.558
BOW-CN	100	2015	-	0.122
CCL VGGM [8]	1024	2016	0.386	-
Mixed Diff+CCL [8]	1024	2016	0.455	-
HDC+Contrastive [23]	384	2017	0.575	-
GSTE [20]	1024	2017	0.724	0.594

great application prospect. Compared with the classic person Re-identification problem, vehicle ReIdentification is more challenging since it faces the enormous inter-class similarity and intra-class variances presented by massive vehicles of the same model types and the shooting situation variations across multiple cameras. For example, the subtle differences between similar vehicles are even challengeable for human beings. Fortunately, some special marks such as tissue box, pendant, annual inspection marks, etc provide characteristics clues for efficient discrimination.

The deep metric learning has been widely adopted for Vehicle ReID tasks. The objective of the deep network is to learn a deep embedding where the samples of the same vehicle ID are constrained in a local space, such that the samples of different vehicle ID are farther away than ones of the same vehicle ID. Such feature distribution is pretty desirable for nearest neighbor retrieval. As such, the retrieval method is the main solution for re-identification. As mentioned above, the granularity of vehicle ReID is finer than person ReID. Consequently, the inter-class feature distribution requires more structured prior knowledge to represent such subtle differences. This motivated some recent attempts which incorporate intra-class variance into the feature representation, such as group sensitive triplet embedding in [20].

#### 4 STANDARDIZATION OF DEEP FEATURE DESCRIPTOR

In the context of video big data, to further ensure interoperability in deep learning based video analysis, a standard that focuses primarily on defining the syntax of compressed deep feature descriptors is essential. This section clarifies the issues to be solved in the standardization process and how they might be pragmatically approached. We believe

Table 4

Performance summarization of some representative Person ReID methods on the benchmarks

Methods	Dims	Year	VIPeR	CUHK1	PRID	GRID
Cov-of-Cov [24]	16828	2014	33.9	40.9	47	16.6
LOMO [25]	26960	2015	40	-	15.3	16.6
GOLD [26]	1169	2015	27.1	35.3	40.5	10.9
2AvgP [27]	952	2015	28.8	36.1	44.7	12.9
GOG-RGB [17]	7567	2016	42.3	55.8	63.6	22.8
NFST [28]	5138	2016	51.2	69	-	-
SCSP [16]	120	2016	53.5	24.2	-	-
SSDAL [15]	105	2016	43.5	-	20.1	19.1
TMA [29]	100	2016	39.9	-	54.2	-
P2S [30]	800	2017	-	77.3	-	-
Spindle [18]	256	2017	53.8	79.9	67	-

Table 5

A summary of efficiency and accuracy comparisons between recent remarkable works on face recognition

Methods	Year	Dims	LWF	YTF
DeepFace [31]	2014	4096	97.35	91.4
Learning Face [32]	2014	10575	97.73	92.2
MDML-DCPs [33]	2015	1024	98.95	97.3
FaceNet [34]	2015	128	99.63	95.12
Deep embedding [35]	2015	128	99.13	-
Multimodal deep face [36]	2015	9000	98.43	-
Center Loss [19]	2016	512	99.28	94.9
Large-margin softmax [37]	2016	512	98.71	-
SphereFace [6]	2017	512	99.42	95
Neural Aggregation [38]	2017	128	-	95.72

that such AI oriented standard could represent a sea change in the future smart city applications.

#### 4.1 Compact Deep Feature for Video Analysis

As introduced in Section 3, the features extracted by deep neural networks are gradually replacing the handcrafted features in many visual intelligence analysis. Due to millions of parameters lying in the deep network, as well as a series of non-linear mappings, the deep network can present high discrimination capability with pretty lower memory costs compared with the handcrafted features. Moreover, when massive training data is available, the involvement of end-to-end learning scheme would further sharpen the feature discrimination ability. Here, we investigate the performance and feature compactness of the recent remarkable works in four typical analysis tasks in city surveillance.

Although different network structures are employed in different analysis tasks, we find that the features can be uniformly represented without significantly sacrificing the analysis accuracy. Table 2 lists the video retrieval performance of the deep learning features with off-the-shelf CNN model reported in CDVA benchmarks. From the perspective of the performance and feature compactness, the deep learning feature shows the competitive performance. With the advance of network structure and training scheme, the performance of deep features have also been dramatically improved, being state-of-the-art on several image retrieval benchmarks, such as Holiday, Oxford5K and Paris. The person/vehicle ReID tasks also benefit a lot from the success of deep networks, as shown in Tables 3&4. In particular, we observe there is a trend that the recent methods employ features with much lower dimensions to produce the

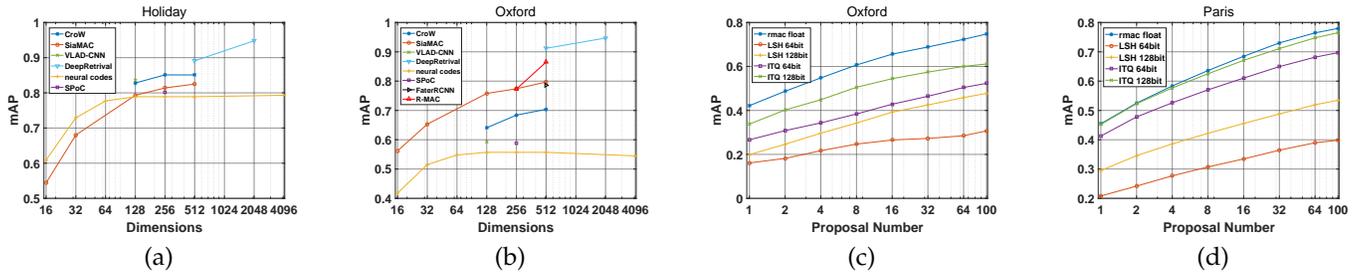


Figure 3. Performance variations with the increase of dimensions and object proposals for recent works on image retrieval benchmarks.

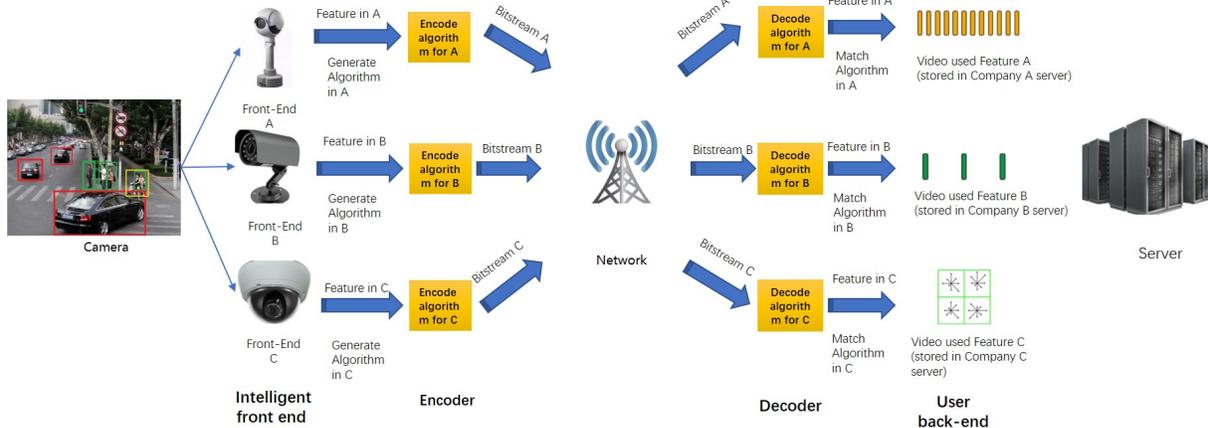


Figure 4. Illustration of traditional feature transmission framework. Different feature bitstreams originated from different front end cameras have their own organization and syntax, such that the corresponding encoding and decoding algorithms should be performed at local and server ends, respectively.

representation of an object, with the help of more powerful networks and well-defined loss functions. Similar trends can also be found in the recent efforts for face recognition, as shown in Table 5. In addition, the performance improvements originate from the introducing of object proposal in image retrieval has been witnessed in several benchmarks such as Oxford and Paris. The performance gains originate from the better localization ability provided by region proposal as shown in Fig. 3.4 (c)&(d). However, the obvious demerit of region-based methods is the relatively high feature dimension due to the representation requirements of multiple targets in query or reference images. As such, the region proposal mechanism would be practically applicable with the representations of compact feature. In Fig. 3.4 (a)&(b), the performance variations against the dimension are illustrated. All these experimental results demonstrate that these tasks can be successfully achieved with identical or similar feature dimension. For example, when the feature size reaches 512 dimensions, in most of the face recognition cases, competitive results can be obtained. Obviously, similar phenomenon can also be found when the dimension reaches 512 in CDVA, 512 in Person ReID, and 1024 in vehicle ReID. In a word, a converging point can be feasibly attained from the perspectives of feature dimensions, proposal numbers, network structure. Another observation is that the performance variation along with the augment of proposals can also arrive at a saturation

eventually as shown in Fig. 3.4 (c)&(d). Such observations provide useful evidence for the further standardization of deep learning features, as discussed in Section 4.2. In the future, it is also expected that more compact and discriminative feature representations will emerge due to the advance of the network architectures and optimization strategies.

## 4.2 Toward Standardization of Deep Features

It is apparent that the compact deep feature possesses many favorable properties for the applications of the smart city. However, the explosion of the deep learning models is also creating many challenging research problems. In particular, it is worth noting that the feature coding differs with traditional video coding in that an end-to-end feature coding pipeline involves both feature extraction and compression. In other words, for video coding the source visual signals are established and available, i.e. the pixel values. By contrast, in feature coding, different deep learning models would create dramatically different features for the subsequent compression process. Therefore, a complete and exhaustive standard that can fully ensure the interoperability typically specifies the standardization of both feature extraction and compression. As such, any bitstreams that conform to such standard can be meaningfully compared.

Such standardization requires the deterministic deep network model and parameters. Nevertheless, the recent research achievements of deep learning emerge in endlessly,

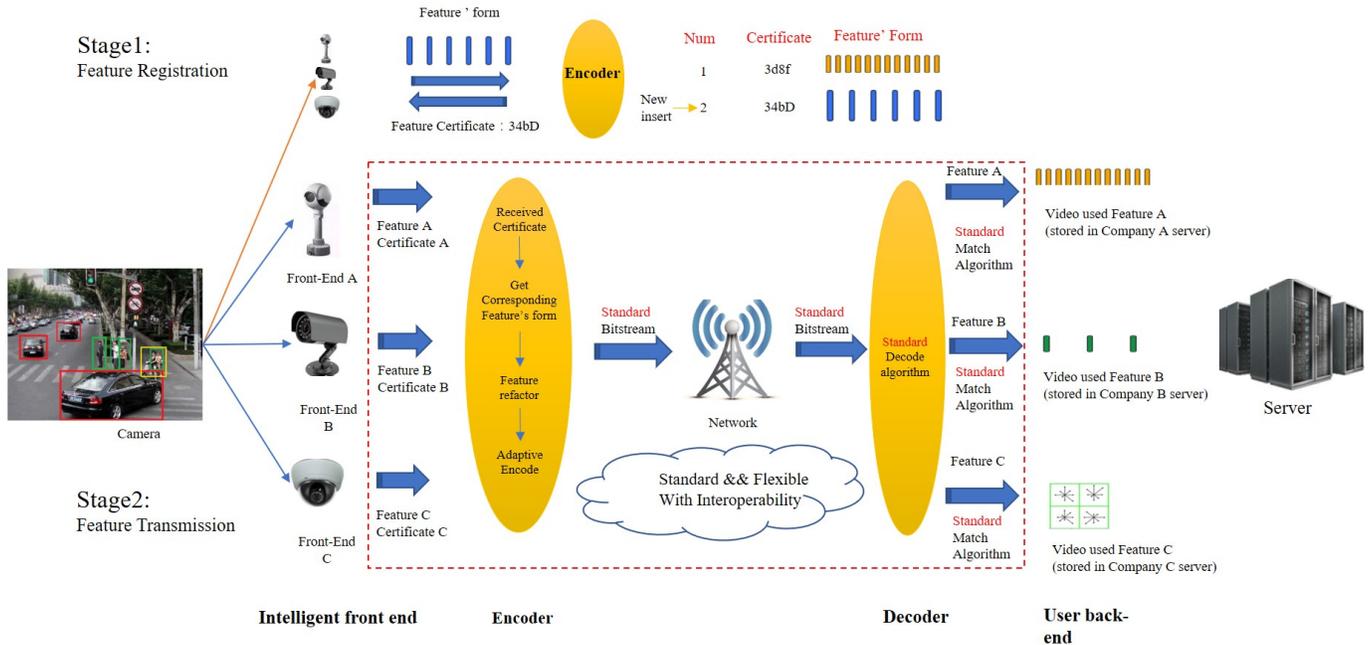


Figure 5. Illustration of the interoperability-enabled feature transmission framework. First, new features are required to register its feature organization and syntax in the encoder and obtain the corresponding certificate. Second, features generated at the front end will be reorganized according to its certificate registered in encoder, and then encoded in standard bitstream syntax. Therefore, The server end can leverage standard algorithm to decode the received features

and moreover, there is a lack of the generic deep model that can be applied to a broad of tasks in video surveillance. Therefore, the standardization of deep learning model is not ready for prime time. Here, we propose the concept of semi-interoperability for feature coding, which only standardizes the feature compression. In other words, only the pipeline from raw features to the compressed bitstream is taken into consideration, and the final syntax that specifies the compact deep features is standardized. The raw feature extraction process is left open for future exploration. Such strategy is based on the key observation that the raw features for these tasks can be uniformly represented, as demonstrated in Section 4.1. As such, the increasing demand for the interoperability in smart city and the explosion of deep learning techniques can be well balanced.

The semi-interoperability based standardization strategy is dual to the video coding standard where only the decoder is standardized. The decoder conforming to the standard can only correctly recover the features, but does not account for the explanation of the features as the deep learning model is not specified. Therefore, such strategy only ensures that any deep learning feature bitstreams from the same deep learning model conforming to such standard can be meaningfully matched after decoding. In other words, it does not fully support the interoperability and bitstreams conforming to such standard may convey different information. On the other hand, the advantage lies in that in the future any effective deep learning models can seamlessly collaborate with this standard, such that the standard can be kept with long-lasting vitality. Moreover, though there are multiple tasks in video surveillance and each task corresponds to the specific deep learning model, as long as the final generalized bitstreams from these models conform to

the standard, they can be successfully decoded by a unified decoder. Here, the traditional feature compression and standardized feature compression frameworks are shown in Figs. 4&5. It is observed that the bit-streams from different ends can be uniformly represented and transmitted, such that a unified decoder can be used to decode such bitstream to enable the semi-interoperability.

Regarding feature compression, the high redundancy of deep learning features in video sequences needs to be removed. In particular, many video coding technologies can be analogously transferred to feature codings, such as inter prediction, intra prediction and rate distortion optimization. In addition, since the basic role of video surveillance is to analyze and explain the object behaviors, and in many occasions within a video frame there are multiple objects, it is natural to extend the frame level feature extraction and compression to the object level based on object proposals. For example, real-time object detectors such as YOLO [39] can be adopted to localize the target objects such as persons, vehicles, heads or other objects of interest, then the regions of interest will be feed into the corresponding networks designed for specific tasks to obtain the feature representation. This also requires the non-local intra prediction to remove the redundancy from different objects within a frame. As such, how these redundancies can be removed and how the final bitstream is composed of should be further investigated in the standardization exploration.

It is also anticipated that in the future the deep learning models are developed to maturation and generic as well as dynamic feature representations can be learned from surveillance videos. At that stage, there may emerge a unified deep learning model that can be standardized to achieve the full interoperability. Generally speaking, such

deep learning model can not only deal with the various video surveillance tasks, but also enjoy the properties such as lightweight and friendly for implementation. Overall, the message we are trying to send here is not that the standardization of deep model for feature extraction is abandoned. Rather, we hope to make the point that at the current stage, there are flexible and practical alternative solutions for the standardization that can be deployed.

## 5 OUTLOOK

We have discussed the practical issues and envisioned the future standardization of deep learning features in the context of large-scale video management in the smart city. Rather than exhaustively establishing the whole feature representation process including both extraction and compression, we have emphasized on the great potentials of standardizing the bitstream syntax of the compressed features. Such strategy is significantly different from the previous MPEG-7 visual standards such as CDVS and CDVA, and the deep learning models are not required to be specified to conform to the standard, which further enhances the flexibilities in the proliferation of deep learning technologies. In the future, it is expected that such AI oriented feature coding standard can play important roles in the establishment of the visual system of the city brain, and impact the new development of future AI technologies.

## REFERENCES

- [1] Ian H Witten, Eibe Frank, Mark A Hall, and Christopher J Pal, *Data Mining: Practical machine learning tools and techniques*, Morgan Kaufmann, 2016.
- [2] Zisserman Andrew Simonyan, Karen, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [3] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [5] Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten, "Densely connected convolutional networks," *arXiv preprint arXiv:1608.06993*, 2016.
- [6] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song, "Sphereface: Deep hypersphere embedding for face recognition," *arXiv preprint arXiv:1704.08063*, 2017.
- [7] Tong Xiao, Hongsheng Li, Wanli Ouyang, and Xiaogang Wang, "Learning deep feature representations with domain guided dropout for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1249–1258.
- [8] Hongye Liu, Yonghong Tian, Yaowei Yang, Lu Pang, and Tiejun Huang, "Deep relative distance learning: Tell the difference between similar vehicles," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2167–2175.
- [9] Ling-Yu Duan, Vijay Chandrasekhar, Jie Chen, Jie Lin, Zhe Wang, Tiejun Huang, Bernd Girod, and Wen Gao, "Overview of the MPEG-CDVS standard," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 179–194, 2016.
- [10] "Call for Proposals for Compact Descriptors for Video Analysis (CDVA)-Search and Retrieval," *ISO/IEC JTC1/SC29/WG11/N15339*, Warsaw, Jun, 2015.
- [11] Jie Lin, Ling-Yu Duan, Shiqi Wang, Yan Bai, Yihang Lou, Vijay Chandrasekhar, Tiejun Huang, Alex Kot, and Wen Gao, "Hnnp: Compact deep invariant representations for video matching, localization, and retrieval," *IEEE Transactions on Multimedia*, vol. 19, no. 9, pp. 1968–1983, 2017.
- [12] Lempitsky Victor Babenko, Artem, "Aggregating local deep features for image retrieval," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1269–1277.
- [13] Giorgos Tolias, Ronan Sifre, and Hervé Jégou, "Particular object retrieval with integral max-pooling of cnn activations," *arXiv preprint arXiv:1511.05879*, 2015.
- [14] Albert Gordo, Jon Almazán, Jerome Revaud, and Diane Larlus, "Deep image retrieval: Learning global representations for image search," in *European Conference on Computer Vision*. Springer, 2016, pp. 241–257.
- [15] Chi Su, Shiliang Zhang, Junliang Xing, Wen Gao, and Qi Tian, "Deep attributes driven multi-camera person re-identification," in *European Conference on Computer Vision*. Springer, 2016, pp. 475–491.
- [16] Dapeng Chen, Zejian Yuan, Badong Chen, and Nanning Zheng, "Similarity learning with spatial constraints for person re-identification," in *Computer Vision and Pattern Recognition*, 2016, pp. 1268–1277.
- [17] Tetsu Matsukawa, Takahiro Okabe, Einoshin Suzuki, and Yoichi Sato, "Hierarchical gaussian descriptor for person re-identification," in *Computer Vision and Pattern Recognition*, 2016, pp. 1363–1372.
- [18] Haiyu Zhao, Maoqing Tian, Shuyang Sun, Jing Shao, Junjie Yan, Shuai Yi, Xiaogang Wang, and Xiaoou Tang, "Spindle net: Person re-identification with human body region guided feature decomposition and fusion," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1077–1085.
- [19] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao, "A discriminative feature learning approach for deep face recognition," in *European Conference on Computer Vision*. Springer, 2016, pp. 499–515.
- [20] Yan Bai, Feng Gao, Yihang Lou, Shiqi Wang, Tiejun Huang, and Ling-Yu Duan, "Incorporating intra-class variance to fine-grained visual recognition," *arXiv preprint arXiv:1703.00196*, 2017.
- [21] Massimo Balestri, Miroslaw Bober, and Werner Bailer, "Cdva experimentation model (cxm) 0.2," *ISO/IEC JTC1/SC29/WG11/N16274*, Geneva, May, 2016.
- [22] Jiang Wang, Yang Song, Thomas Leung, Chuck Rosenberg, Jingbin Wang, James Philbin, Bo Chen, and Ying Wu, "Learning fine-grained image similarity with deep ranking," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1386–1393.
- [23] Yuhui Yuan, Kuiyuan Yang, and Chao Zhang, "Hard-aware deeply cascaded embedding," 2016.
- [24] Giuseppe Serra, Costantino Grana, Marco Manfredi, and Rita Cucchiara, "Covariance of covariance features for image classification," in *International Conference on Multimedia Retrieval*, 2014, p. 411.
- [25] Shengcai Liao, Yang Hu, Xiangyu Zhu, and Stan Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *Computer Vision and Pattern Recognition*, 2015, pp. 2197–2206.
- [26] Giuseppe Serra, Costantino Grana, Marco Manfredi, and Rita Cucchiara, "Gold: Gaussians of local descriptors for image representation," *Computer Vision Image Understanding*, vol. 134, pp. 22–32, 2015.
- [27] Jo Carreira, Caseiro Rui, Jorge Batista, and Cristian Sminchisescu, "Free-form region description with second-order pooling," *IEEE Transactions on Pattern Analysis Machine Intelligence*, vol. 37, no. 6, pp. 1177–1189, 2015.
- [28] Li Zhang, Tao Xiang, and Shaogang Gong, "Learning a discriminative null space for person re-identification," in *Computer Vision and Pattern Recognition*, 2016, pp. 1239–1248.
- [29] Niki Martinel, Abir Das, Christian Micheloni, and Amit K. Roy-Chowdhury, "Temporal model adaptation for person re-identification," in *European Conference on Computer Vision*, 2016, pp. 858–877.
- [30] Sanping Zhou, Jinjun Wang, Jiayun Wang, Yihong Gong, and Nanning Zheng, "Point to set similarity based deep feature learning for person re-identification," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [31] Yaniv Taigman, Ming Yang, Marc Aurelio Ranzato, and Lior Wolf, "Deepface: Closing the gap to human-level performance in face

- verification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1701–1708.
- [32] Yi Sun, Xiaogang Wang, and Xiaoou Tang, "Deep learning face representation from predicting 10,000 classes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1891–1898.
- [33] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, et al., "Deep face recognition," in *BMVC*, 2015, vol. 1, p. 6.
- [34] Florian Schroff, Dmitry Kalenichenko, and James Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 815–823.
- [35] Jingtuo Liu, Yafeng Deng, Tao Bai, Zhengping Wei, and Chang Huang, "Targeting ultimate accuracy: Face recognition via deep embedding," *arXiv preprint arXiv:1506.07310*, 2015.
- [36] Changxing Ding and Dacheng Tao, "Robust face recognition via multimodal deep face representation," *IEEE Transactions on Multimedia*, vol. 17, no. 11, pp. 2049–2058, 2015.
- [37] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang, "Large-margin softmax loss for convolutional neural networks," in *ICML*, 2016, pp. 507–516.
- [38] Jiaolong Yang, Peiran Ren, Dong Chen, Fang Wen, Hongdong Li, and Gang Hua, "Neural aggregation network for video face recognition," *arXiv preprint arXiv:1603.05474*, 2016.
- [39] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali 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.