

# Deep Video Representation Learning: a Survey

Elham Ravanbakhsh<sup>1</sup>, Yongqing Liang<sup>2</sup>, J. Ramanujam<sup>1</sup>, Xin Li<sup>3,2\*</sup>

<sup>1</sup>Division of Electrical & Computer Engineering and Center for Computation & Technology, Louisiana State University, Baton Rouge, 70803, LA, USA.

<sup>2</sup>Department of Computer Science and Engineering, Texas A&M University, College Station, 77843, TX, USA.

<sup>3</sup>Section of Visual Computing and Interactive Media, Texas A&M University, College Station, 77843, TX, USA.

\*Corresponding author(s). E-mail(s): [xinli@tamu.edu](mailto:xinli@tamu.edu);  
Contributing authors: [eravan1@lsu.edu](mailto:eravan1@lsu.edu); [lyq@tamu.edu](mailto:lyq@tamu.edu); [jxr@cct.lsu.edu](mailto:jxr@cct.lsu.edu);

## Abstract

This paper provides a review on *representation learning for videos*. We classify recent spatio-temporal feature learning methods for sequential visual data and compare their pros and cons for general video analysis. Building effective features for videos is a fundamental problem in computer vision tasks involving video analysis and understanding. Existing features can be generally categorized into spatial and temporal features. Their effectiveness under variations of illumination, occlusion, view and background are discussed. Finally, we discuss the remaining challenges in existing deep video representation learning studies.

**Keywords:** Video Representation Learning, Feature Modeling, Video Feature Extraction, Feature Learning.

## 1 Introduction

The enormous influence of media and social networking has led to an avalanche of videos uploaded on the internet every day. To effectively analyze and use the uploaded video data, it is important to construct feature representations for videos. Unlike the analysis and understanding of images, the manual modeling of video features is often a

laborious task. Therefore, there is a need for techniques that can automatically extract compact yet descriptive features.

With the recent advances in artificial intelligence and computer vision, deep neural networks have achieved significant success in feature modeling. These techniques have led to a great breakthrough in practical video analysis applications such as tracking [91, 177], action recognition [187], action prediction [90], and person re-identification [51]. To design a deep learning pipeline for these applications, extracting video features is often the first step and it plays a critical role in subsequent video processing or analysis. Developing deep learning pipelines to extract effective features for a given video is referred to as *deep video representation learning*.

The characteristics of videos are often encoded by *spatial features* and *temporal features*. Spatial features encode geometric structures, spatial contents, or positional information in image frames; whereas, temporal features capture the movements, deformation, and various relations between frames in the time domain. Depending on the target applications, an algorithm should be able to extract either spatial or temporal features, preferably, both. Some applications also require the decoupling of spatial and temporal information from the extracted features so that some specific characteristics can be more effectively modeled.

Learning robust representation for videos faces several *major challenges* such as

- occlusion: objects of interests might be partly occluded;
- illumination: videos might be taken under various lighting conditions or/and from changing view angles;
- view and background variations: foreground objects and background scenes can be moving.

Therefore, we evaluate the performance of representation learning algorithms using *robustness* and *accuracy* under these scenarios. Robustness of different algorithms is evaluated under these four challenges: occlusion, view, illumination, and background change. As for their accuracy, since different applications use different metrics, we adopt accuracy metrics from representative tasks of action recognition and video segmentation, in which more expressive features generally lead to better accuracy.

**Comparison with Existing Surveys.** Representation learning from images is a classic problem in computer vision, and it has been widely studied to facilitate various image analysis and understanding tasks. Many survey papers have been published to address this problem. But most of these representation learning studies focused on features of static images [63, 112, 114, 124]. As summarized in Table 1, while multiple components of these studies are closely related to video representation, a systematic survey on video features is missing. Some recent surveys [43, 61, 86] discussed video representation learning, but most of these have focused on a specific type of learning or method. A few other survey papers discussed video processing tasks that involve video representation learning [25, 57, 64, 144, 167]; however, their focuses were mostly on discussing how the developed pipelines perform on the targeted task(s). There is a lack of a survey of representation learning in a general setting and one that investigates the role of each feature on its embedding regardless of its specific task. We believe our survey provides an insightful analysis on how to construct effective

video features/representations, particularly when confronting with aforementioned challenges.

1. We provide a comprehensive survey of deep video representation learning.
2. We compare different types of representation learning algorithms in terms of accuracy and robustness in various practical challenging scenes, and provide some observations/suggestions in adopting suitable features for different video processing and analysis tasks.

**Table 1** Current surveys related to image and video feature learning. Unlike existing surveys that studied representation learning for specific targeting tasks, our survey discusses pros and cons of different features for general scope.

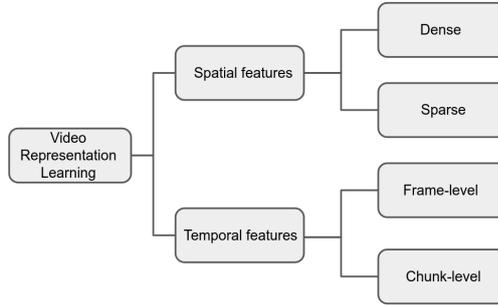
Reference	Year	Image	Video	General scope	Application
[124]	2018	✓	-	×	Visual-based localization
[112]	2021	✓	-	×	Image matching
[114]	2021	✓	-	×	Semantic and instance segmentation
[63]	2021	✓	-	×	Content-based image retrieval
[43]	2017	✓	✓	×	Representation learning on graphs
[86]	2018	✓	✓	×	Multi modal learning
[61]	2020	✓	✓	×	Self supervised feature learning
[141]	2021	✓	✓	×	2D and 3D pose estimation
[144]	2021	✓	✓	×	Multi person pose estimation
[25]	2020	-	✓	×	Soccer video analysis
[57]	2021	-	✓	×	Multi-view video summarization
[64]	2022	-	✓	×	Event detection in surveillance videos
[167]	2022	-	✓	×	Pedestrian attribute recognition
<b>Ours</b>	<b>2023</b>	✓	✓	✓	<b>Action recognition + video segmentation</b>

**Organization.** We present the classification of deep features for videos in Section 2, and then compare these features in action recognition and video segmentation application in Section 3. We conclude the paper by discussing remaining challenges and future directions in Section 4.

## 2 Classification of Deep Video Features

Two main aspects of video data are often considered in video processing and analysis tasks: (1) how to encode spatial structures or contents in visual data; and (2) how to model temporal coherency or changes among frames. Appearance information and geometry structure of the scene or objects are considered as **spatial information**. In representation learning, capturing spatial relation is important in understanding the visual concept of the video. Based on spatial information extracted, we generally divide features into **dense** features and **sparse** features. **Spatially dense** features are contextual data often defined using pixel intensities of the input. Typical and widely used dense features can come from RGB images and their variants, such as RGBD images. **Spatially sparse** features are often defined on a smaller set of entities such as keypoints, subpatches, or other graph structures. Widely used sparse features include

those defined on divided patches or structural graphs (e.g., media axes for general shapes, skeletal structures of humans/animals).



**Fig. 1** Classification of deep video representation learning schemes.

**Table 2** Pros and cons of different types of features. Scene noise includes occlusion, illumination, background and viewpoint variations.

Features	Pros	Cons
Dense	Good in capturing appearance information	Sensitive to scene noise High intra-class variations High computational cost
Sparse	Robust against background and illumination change Low intra-class variations Low computational cost	Weak in capturing appearance information
Frame-level	Low computational cost	Weak in co-occurrence representation learning
Chunk-level	Good in co-occurrence representation learning	High computational cost

The second aspect is **temporal information** which plays an important role in video representation learning, and is a key difference between image features and video features. To effectively understand the concept of a video, temporal information or temporal coherence across different frames plays a critical role. Based on how the temporal information is modeled, we divide temporal features into two categories: **frame-level** and **chunk-level** features. The former extracts features from each frame and constructs a sequence of signatures; while the later one encodes spatio-temporal features of a chunk into one signature.

We illustrate our classification in Fig. 1. In the following subsections 2.1 to 2.4, we classify video features based on how their spatial and temporal information is modeled, and discuss their pros and cons, as summarized in Table 2. We will discuss how different designs affect the features' general robustness under different scenarios.

In terms of feature accuracy, it is more application-dependent, and will be discussed in the application section.

## 2.1 Spatially dense features

Spatially dense features contain rich information mostly defined by using direct pixel intensities of the input in a structured order. RGB video is the most common data type in dense features. It contains contextual data about appearance, objects and background. People often use a 2D Convolutional Neural Network (CNN) as a *standard architecture* for extracting spatial information from dense features.

**CNN-based Standard Architectures.** With the emergence of deep learning, CNNs have become the most common method for feature modeling due to the strong modeling capability and superior performance of deep learning-based methods. Various CNN architectures including VGGNet [140] and ResNet-50 [46] have been used for spatial representation learning [34, 97, 200]. These models provide a high-level spatial representation of video frames. Some approaches also use object detection or segmentation algorithms to extract local regions of interest in a frame to better exploit the correlation of different regions and reduce the chance of encoding redundant information [5]. For example, in [90], an interaction module is proposed to model interactions between a (human) object and its surroundings. Mask-RCNN [47] is used for semantic segmentation, then these masked features are given to a 2-layer convolutional network for further feature learning.

**Extra modules for better robustness.** Effective video features should be robust against occlusion, view and background change. Therefore, built upon the above *standard architectures*, people also add extra modules to improve feature robustness. Recent approaches adopted three general types of additional modules: *part information*, *additional input information*, and *attention mechanism*. In the following, we will elaborate that given a *standard architecture*  $X$ , adding one or multiple modules onto  $X$  could improve the pipeline’s robustness under different scenarios.

### *Part information*

Typical video representation learning methods do not take into account the effect of partial occlusion. But with partial occlusion in the video, learned features are often corrupted due to the inclusion of irrelevant surrounding objects, and consequently, cause dramatic performance degradation. A conventional strategy is to train ensemble models for various occlusion patterns [15, 118, 149], and construct a part pool that covers different scales of object parts, and then automatically choose important parts to overcome occlusions. A main limitation of this strategy is that it is more expensive in both training and testing phases. Another strategy is to integrate a set of part/occlusion-specific detectors in a joint framework that learns partial occlusions [119, 203]. In [203], a set of part detectors are applied and the correlation among them is exploited to improve their performance. Each part detector generates a score for a candidate region and the final score is computed based on the average among some top scores. The main issue with such part-specific detectors is that they are not able to cover all the occlusion patterns comprehensively and need to be designed based on

some pre-assumptions [183]. For example, for pedestrian detection, these part detectors are designed with the prior knowledge that pedestrians are usually occluded from the bottom, left, and right. However, in practice, occlusion patterns can be irregular, which affects the feature’s performance.

#### ***Additional Information***

Some studies use additional information to compensate for the view and the illumination variance of dense features. For example, RGB-D cameras are often used to include depth information which makes the representation less sensitive to illumination and view variations [52, 100]. Thermal images together with color images are also used to improve robustness when suitable light source is not available [71]. Several other works use non-vision information as a complementary input. For example, in [40], audio signals and in [105] signals from wearable sensors are added to improve robustness of dense features.

#### ***Attention Mechanism***

In most scenarios, key objects/regions are just part of the whole spatial image. Being able to use local spatial attention to guide the system to focus on important (foreground) object and ignore irrelevant (background) noise is desirable. Many recent studies developed attention mechanisms to help feature learning models concentrate on important regions in the spatial dimension. For example, in [183], an attentive spatial pooling is used instead of max-pooling which computes similarity scores between features to compute attention vectors in the spatial dimension. This method allows model to be more attentive on region of interests in image level. In [88], a self-attention mechanism is adopted to generate a context-aware feature map. Then the similarity between context-aware feature maps and a set of learnable part prototypes are calculated and used as spatial attention maps. To identify the object of interest and focus on visible parts, some models use spatially sparse features in spatial attention networks to construct a robust representation. In [178, 198], sparse features aid RGB frames to estimate the attention map and visibility scores to handle various occlusion patterns and generate view-invariant representation.

#### ***Summary of Extra Modules***

Although dense features can encode rich contextual information, their representation can sometimes be sensitive to occlusion, view, illumination, and background variance. Given a standard architecture  $X$  that extracts dense features, various recent studies exploited adding extra modules onto  $X$  to enhance the robustness of spatial feature modeling. Table 3 summarizes these strategy discussed above.  $X$  + part information is often effective in enhancing the model’s robustness against partial occlusion.  $X$  + additional input information (e.g., depth info) can help the model enhance its performance under illumination and view changes. Many recent studies incorporate attention mechanism into the standard architecture, such an attention strategy often help differentiate foreground objects of interest and the background noise, and consequently, helps the model perform better towards view, occlusion, and background variance.

**Table 3** Robustness of different models using Dense features. A standard architecture is called X. Its limitations against robustness challenges can be addressed/alleviated by adding extra modules.

Model	View	Occlusion	Background	Illumination
X	-	-	-	-
X + part information	-	✓	-	-
X + additional information	✓	-	-	✓
X + attention	✓	✓	✓	-

## 2.2 Spatially sparse features

Spatially sparse features represent information with a sparse set of feature points that are able to describe the geometry of original video frames. There are several advantages of using sparse features. First, they describe frames with a sparse set of features that leads to a low computational cost. Second, sparse features suffer less intra-class variances compared to dense features and are more robust to the change of conditions such as appearance, illumination, and backgrounds. However, sparse features lack appearance information which in some scenarios is essential for feature modeling. Three types of *standard architectures* are often adapted to extract sparse features: *RNN*, *CNN*, and *GNN/GCN pipelines*.

**RNN-based Standard Architectures.** Due to the recurrent nature of videos, one strategy is to use RNN-based architectures for sparse representation learning. Spatial feature arrangements determine the proper architecture for spatial extraction. If the inherent spatial relations between feature points can be arranged in a sequence of vectors or grids, RNNs are effective for spatial modeling [32, 156].

*Pros and Cons of RNN-based Standard Architectures* RNN-based methods are good in dealing with sequential data, but they are not very effective in spatial modeling. Therefore, their general performance is not as good as CNN models [102]. In several works, RNN-based models take feature maps from CNNs as their input rather than using raw input frames [34, 181].

**CNN-based Standard Architectures.** There are many studies that deployed CNNs for sparse representation learning. However, defining the relations between unstructured feature points is challenging. CNNs take their input in a form of an image. To satisfy the need of CNNs’ input, some researches model the sparse features into multiple 2D pseudo-images [75, 101, 104]. For example, in [101], sparse features are arranged into an image containing the feature points in one dimension and frames in another one, then a 7-layer CNN is applied to extract features. In [164], the frame dynamics are mapped into textured images, then a CNN model is used to extract information. Some works use a group of feature points to construct a hierarchy among features. In [48], the video frame is divided into multiple grids, then a convolutional feature descriptor is run for each cell.

*Pros and Cons of CNN-based Standard Architectures* CNN-based architectures are effective in extracting local features and exploring discriminative patterns in data. However, their major drawback is that they are designed for image-based input and primarily rely on spatial dependencies between the neighboring points. For some sparse

features that contain unstructured design, CNNs cannot perform very well. Additionally, CNNs have trouble with wildly sparse data as they heavily rely on spatial relations of points to learn. In that case, RNNs perform better.

**GNN/GCN Standard Architectures.** Modeling sparse features in a vector or an image may corrupt the spatial relations or add false connections between feature points that their relation is not strong enough. Some of these irregular features are intrinsically structured as a graph and cannot be modeled in a vector, 2D or 3D grid. To address this issue, some researches utilize graph structure, such as graph convolutional neural networks (GCNNs) [76, 80, 187] or GNNs [44] to build feature extraction architectures.

Some people use spatio-temporal GCNs by arranging feature points as an indirect graph with points as nodes and their relations as edges. Nodes within one frame are connected based on the spatial relations of the features and represented by an adjacency matrix in spatial dimension. Nodes in the temporal dimension are related through the relations of corresponding nodes in consecutive frames [187]. A concern with spatio-temporal GCNs is that modeling spatial features just according to natural connection of nodes might lose the potential dependency of disconnected joints. To solve this issue, in [80, 133], a model is introduced to adaptively learn and update the topology of graph for different layers and samples. In [80], a framework is proposed to capture richer dependencies among points and neighbors. They used an encoder-decoder structure to model dependencies between far-apart points. When modeling sparse features, defining connections between different feature nodes is challenging. For example, spatio-temporal GCN models the features as an indirect graph, which may not fully express the direction and position of points. [133] used a directed acyclic graph to represent sparse data. Then a directed graph neural network is used to encode the constructed directed graph, which can propagate the information to adjacent nodes and edges and update their associated information in each layer.

*Pros and cons of GNN/GCN Standard Architectures* GNN/GCNs are helpful in dealing unstructured data that have underlying graph structures and are non-Euclidean. The non-regularity of data structures in most sparse features have led to superior performance of graph neural networks over CNN or RNN architectures. But if sparse features can be formulated as 2D grids or vectors, then graph-based models are not as efficient as CNN models.

**Extra Modules for Better Robustness.** The representation of spatially sparse features is robust to background and illumination change. To overcome other robustness issues, some studies add extra modules on the standard architecture. Standard architecture X could be a CNN, RNN or GNN/GCN model. Generally, two extra modules, *transformation matrix* and *attention*, are often added to enhance feature robustness.

### ***Transformation matrix***

Change of camera view points can change the relative position of feature points. Hence, constructing a view-invariant representation remains still a challenge in sparse feature modeling. To address this problem, several researches developed a transformation method to transform a set of feature points to a standard structure [56, 87, 104]. In

[87], a rotor-based view transformation method is proposed to re-position the original features to a standard frontal system. After transformation, a spatio-temporal model is applied to construct the shape and motion of each part. In [104], a sequence-based transformation is applied on the features to map them to a standard form and make a new view-invariant sequence.

### ***Attention mechanism***

Attention aids model in confronting with partial occlusion and view variations. Similar to dense features, sparse feature modeling in the presence of occlusion is still challenging. Some approaches deploy spatial attention mechanism to focus on the points of interest and predict occluded parts by the help of visible feature points in the adjacent frames. In [51], a spatial attention generator is proposed to predict occluded parts. The generator is an autoencoder that predicts the content of occluded part conditioned on the visible parts of the current frame. Some other studies use heatmaps as attention mechanism to focus on informative feature points. In [35], heatmaps for the visible and occluded part are generated. Then, using the heatmaps occluded parts are predicted along both spatial and temporal dimension. In [142], class activation maps are used as a mask matrix to force network has to learn features from currently inactivated points. To achieve a view-invariant representation, some models propose transferring attention from reference view to arbitrary views. For example, in [60], attention maps are produced to transfer attention from a reference view to arbitrary views. This helps learn effective attention to crucial feature points. In [163], attention directly operates on network parameters rather on input features. This allows spatial interactional contexts to be explicitly captured in a unified way.

*Summary.* Table 4 shows the robustness of standard architecture X and extra modules. X could be RNN, CNN or GNN/GCN architecture and is robust against background and illumination changes. While pros and cons of X depend on the target task, generally RNN is better suited to more sparse data while CNN deals better with denser feature points. Typically if the nature of the structure of data is in a grid format, CNN is more effective in extracting spatial features. On the other hand, if data have a graph scheme, then GNN/GCN performs better. X + transformation matrix is robust against view variations by transforming a set of feature points to a standard view. X + attention is robust against occlusion and view variations by predicting occluded parts from visible points and transfer attention from one view to another, respectively.

**Table 4** Robustness of different models using spatially sparse features. A standard architecture is called X.

model	view	occlusion	background	illumination
X	-	-	✓	✓
X + transformation matrix	✓	-	✓	✓
X + attention	✓	✓	✓	✓

### 2.3 Frame-level features

Frame-level features are a sequence of signatures that each describes a frame individually. Given a video clip, a frame-level feature model often processes spatial information frame by frame and encodes the temporal relationship between frames. People adopt different strategies, such as *Optical flow*, *CNN-based architectures*, *RNN-based architectures*, and *Attention mechanisms* to extract frame level features.

**Optical Flow.** Optical flow is a feature containing motion information of consecutive frames that is useful for describing video dynamics [49]. It is computed by removing the non-moving scene that generates a background invariant representation compared to the original frames. As optical flow is computed frame by frame we categorize it as a frame-level feature. Most studies use optical flow as a temporal data along with their original spatial features to capture movements. Studies in [12, 160, 161], showed that using optical flow and RGB frames achieves a superior performance in modeling videos than only using RGB frames. In [161], optical flow of consecutive 10 frames and RGB are fed to a CNN. The convolutional filters compute derivatives of the optical flow and learn motion dynamics w.r.t the image location. Although to some extent, extracting optical flows as an additional information helps models be more view-independent, robust to occlusion and cluttered background, it often cannot capture relatively long term dependencies. In addition to optical flow, there are other features that contain motion information. For example, in [24], a new motion representation called Potion is proposed that provides all the dynamic of an instance throughout the video clip in one image. Motion cues are represented by different colors which shows the relative time of the frame in the video clip. Using Potion along with RGB frames and optical flows aid model to encode relatively longer dependencies. In [210], in addition to optical flow, motion saliency [18] is calculated from consecutive video frames. Pre-computing motion information including optical flow is time-consuming and storage-demanding. Also, they cannot learn global dependencies among frames.

**CNN-based Architectures.** As CNNs have been widely used in the image tasks, some studies [45, 94, 162] adopted 2D CNNs to model the video clips. However, 2D CNNs need a fusion module to concatenate temporal cues. For example, in [65], late fusion is adopted which fuses information from two different CNNs in the first fully connected layer. In [171], the authors use 2D CNN networks to encode each frame into feature maps then concatenate them as a long vector as the video-level feature descriptor to holistic understanding of the video. However, using a simple fusion or average lacks the learning capability and doesn't contain useful time-related features. In [94], the channels along the time dimension is shifted to improve the performance of temporal modeling with 2D CNNs. In [186], consecutive frames are aggregated with adaptive weights and then fed to convolutional networks.

Recent papers extend the single pass feature encoder to the Siamese structure for time-related feature learning. The Siamese pipeline takes two frames as inputs and uses CNN encoders to extract feature maps. In [17, 54], their feature encoders shared the same weights to extract feature maps from the input images. Then they compare the similarity between the feature maps to build the inter-frame features. In [81, 92, 117], they use two independent encoders to extract features maps. They encode the features

**Table 5** Pros and cons of various architectures used for modeling frame-level features

Models	Pros	Cons
Optical flow	Effective for simple and local dynamics	Ineffective for long and complex dynamics Hard to handle temporal scale variance Expensive computation and storage
2D CNN	Effective for simple and local dynamics	Ineffective for long and complex dynamics hard to handle temporal scale variance
RNN	Effective for complex dynamics	Hard to handle temporal scale variance
Attention	Effective for long and complex dynamics Effectively handle temporal scale variance	Expensive computation

into two parts, one is for similarity comparison, the other stores semantic information of frames. These video features are more useful in the down-stream tasks.

**RNN-based Architectures.** Due to the sequential nature of videos and the ability of memorizing the temporal relations in RNN-based architectures, some studies used these models for encoding motions in a video. Some studies used parallel stream architecture which one stream is responsible for extracting spatial information and the other stream is responsible for extracting temporal data. In [159], a two-stream architecture is proposed where the first stream is responsible for learning spatial dependency and the second for learning the temporal dynamics. The two streams are then aggregated with each other to represent data. As RNNs are not effective in learning long-term temporal dependencies, most models adopt LSTM models. In [34, 200], LSTM is fed with the spatial representation of each frame at each time step. Learning process at each time step is based not only on the observations at that time step, but also on the previous hidden states that provide temporal context for the video. Some other versions of Recurrent-based networks were also explored. For example, in [34], an extended GRU is proposed to model temporal relations by using current and old information. GRU requires less storage and performs faster.

**Attention Mechanism.** As not all frames are informative, several studies used temporal attention to discriminate and select key frames by assigning weights to them. In [13] a scaled dot product attention module is adopted that assigns weights to frame features according to their importance. In [183], all time steps resulting from RNN are combined by an attentive temporal pooling to compute an attentive score in temporal dimension to weight frames based on their goodness. In [197], a non-parametric self and collaborative attention network is proposed to efficiently calculate the correlation weights to align discriminative frames. In [108], a transformer module is adopted that iteratively chooses one frame as query and the rest as key features to compute the temporal attention.

*Summary.* As shown in table 2, while frame-level features have lower computational costs compared to chunk-level features, they decouple spatial and temporal dimensions which leads to lack of co-occurrence representation learning. People adopt different methods for learning frame-level dynamics in a video as shown in table 5. Optical flow and CNN-based architectures assist network to capture short and simple cues. However, they suffer from the lack of memory and therefore they are not suitable

for capturing long term dependencies. Additionally, optical flow is computationally expensive and storage demanding. RNN-based architectures, particularly LSTMs, are able to encode complex dynamics, however they treat each video clip equally and are invariant to inherent temporal diversities. Attention selects informative features and is suited for distinguishing complex tasks.

## 2.4 Chunk-level features

**Table 6** Pros and cons of various dynamics used for modeling chunk-level features

Models	Pros	Cons
3D CNN	effective for coarse level dynamics	ineffective for fine dynamics expensive computation
Attention	effective for fine and coarse level dynamics effectively handle temporal scale variance	expensive computation

While frame-level features separate spatial and temporal modeling completely, chunk-level features extract appearance and dynamics at the same time by creating a hierarchical representations of spatio-temporal data which leads to extracting more subject-related information. These features aggregate frames into one signature, then apply a deep neural network to extract both spatial and temporal information at the same time. People use different strategies to encode chunk-level features: *CNN-based architectures* and *Attention*.

**CNN-based Standard Architectures.** 3D CNN encode a chunk of frames into one signature and has the kernel size of  $s \times s \times d$  which  $s$  and  $d$  refers to the kernel spatial size and the number of frames in one signature, respectively. In [72, 73, 184], a video was split into multiple segments, each of which is fed to a 3D CNN for feature extraction. This 3D CNN starts by focusing on the spatial relations for the first frames and then learns temporal dynamics in the following frames [151]. In [186], two 3D CNN modules are utilized to encode both long and short temporal dependencies by taking chunks with different sizes.

Although 3D CNNs seem like a natural algorithm for modeling video, they have some drawbacks. First, they require a large number of parameters due to adding temporal dimension which leads to leveraging shallow architectures. In [12], a new model was introduced which used a 3D CNN with pre-trained Inception-V1 as a backbone. To bootstrap from the pre-trained ImageNet models [29], the weights of the 2D kernels are repeated along the time dimension. However, the fixed geometric structures of 3D convolution limits the learning capacity of the 3D networks.

Second, chunk-level features learn temporal abstraction of high-level semantics directly from videos, however they are not suited for specific tasks that require granularity in time. In some applications, one may need to precisely predict dynamics in each frame. In this case, chunk level features loose granularity and cannot perform well. For instance, the temporal length of an input video is decreased by a factor of 8 in layers from conv1a to conv5b in C3D architecture [151]. This conforms that local information are lost passing multiple convolutions. To address this issue, in [137], a convolution

and deconvolution approach is proposed which downsamples and upsamples in space and time, simultaneously.

Last but not least, while 3D CNNs are well suited for capturing global coarse motions, they are limited in modeling finer temporal relations in a local spatio-temporal window. To address this, people use attention mechanism to exploit the temporal discriminative information in a chunk of frames.

**Attention Mechanism.** Chunk-level features usually learn the temporal domain by equally treating the consecutive frames in a chunk, while different frames might convey different contributions to the related task. Similar to the frame-level features, there are several works that explore temporal attention in chunk-level features [69, 78]. In [78], attention is learned at channel level by modeling the differences among the temporal channels in 3D CNNs.

*Summary.* Chunk level features are suitable for capturing structural co-occurrence. They connect both spatial and temporal domains and are suitable for learning dynamics that differ in the order of their micro-movements. However, compared with frame-level features they are less suitable for tasks that require fine granularity in time. As shown in table 6, 3D CNNs can capture coarse level temporal dependencies in one signature. But they require a large number of parameters and cannot encode fine details in a local window. Attention is used to exploit informative features in both fine and coarse level and handle temporal scale variance. However, it adds extra weights to the model and is computationally expensive.

### 3 Applying Deep Features in Video Analysis Tasks

After discussing different feature modeling strategies, we compare their usages on different applications. We use two applications, namely, *action recognition* and *video object segmentation*, to analyze these features' behaviors under different circumstances. *Action recognition* aims to recognize specific actions happened in a video and output one (or several, if there are multiple actions) global labels. *Video object segmentation* aims to identify and segment in pixel-level the objects of interest from background in every frame of a video.

In both of these two applications, free-form deformations of objects are common in both spatial and temporal dimensions. And sometimes, to achieve better real-timer efficiency, temporal (or spatial) sampling is intentionally made sparse. Therefore, extracting reliable and powerful features plays a critical role and often directly dictates the final performance.

#### 3.1 Action Recognition

A key challenge in action recognition is how to learn a feature that captures the relevant spatial and motion cues with a descriptive yet compact representation. In the spatial dimension, some studies adopted dense features while some others used sparse features. Likewise, in the temporal dimension, some adopted frame-level features, while some others used chunk-level features. There is a lack of study that discusses the pros and cons of adopting these features under different circumstances. Hence, here we

discuss and analyze each type of features, their challenges, limitations, and possible solutions as shown in table 7.



**Fig. 2** Dense (RGB frames) and sparse (skeleton keypoints) features in Action Recognition. Dense features may include background information; while sparse features encode mainly essential object structure. The usefulness of background varies: it can either distract (e.g., in dancing images) or assist (e.g., activities in the right two columns) recognition. RGB images in the upper left four columns are from [123]; we put them in a sequence of frames. Skeleton keypoints in the lower left four columns are from [111]. Images in the right two columns are from UCF101 dataset [143].

**Table 7** Pros and cons of spatial and temporal features in action recognition. There are some abbreviations in the table: ap., info., bg, rep. means appearance, information, background, and representation respectively.

Features	Pros	Cons	Solution	References
Dense	contains ap. info.	bg noise/redundant info. rep. bias	additional info. attention calibrated data attention	[166, 201] [42, 85]
Sparse	robust to view/bg change low computation cost	low reliability lack of scalability	additional info. heatmap rep.	[11, 33] [24, 185]
Frame	low computation cost	co-occurrence rep.	message passing	[138, 189]
Chunk	co-occurrence rep.	low computation cost fix-length chunk	disentangle kernels multi scale kernels	[152, 162] [26, 58]

**Dense Features.** Due to the adaptability and availability of RGB video frames, such dense pixel-level representations are widely adopted for action recognition [38, 154, 176].

*Pros.* Dense features contain appearance information which is useful in recognizing actions in different scenes. As CNNs have shown their strong ability in capturing dense spatial features, majority of studies use either 2D or 3D CNNs to extract spatial semantics in video frames. Thanks to CNNs, modeling Dense features are straightforward compared to sparse features, however they have their task-related limitations.

*Cons.* There are several challenges in using Dense features. First, they may contain background noise and redundant information which undermines robustness of action representation learning. Another limitation of dense features is “representation bias”

which refers to recognizing actions based on object/background detection. The network may predict correct results based on scene context rather than human action. Some actions might be easier to be predicted using the background and context, like a ‘basketball shoot’ vs a ‘throw’; some others might require paying close attention to objects being interacted by the human, like in the case of ‘drinking from mug’ vs ‘drinking from water bottle’ as shown in Fig. 2. It is noted to mention that representation bias is different from background noise. In background noise, the representation for each class of action differs with the change of the scene, while representation bias means getting help from the discriminative objects in the scene to recognize the action. While some studies in action recognition consider representation bias undesired [23, 85], it may be useful in some scenarios or similar tasks [42]. The reason is that modeling human actions often requires understanding the people and objects around them. For example, recognizing “listening to others” is not possible unless the model knows about the existence of another person in the scene saying something. The main concern with representation bias is that if the dataset is not generalized enough, it can undermine the performance of the model.

*Solutions.* To alleviate the background noise and redundant, several researches have augmented additional visual information to guide network from distraction. Some approaches adopted depth information to overcome the illumination and viewpoint variations [126, 165, 166], while others leverage skeleton data in the form of local attention to assist network in capturing the most representative body postures [110, 201].

Solutions to overcome representation bias include collecting well-calibrated datasets [85], or using an attention mechanism to help the model focus on distinguishable parts of action [132]. Attention networks add a dimension of interpretability by capturing where the network is focusing when modeling actions. In [132], the CNN produces a feature cube for each video input and predicts a softmax over locations and the label classes which determines the probability with which the model believes the corresponding region in the input frame is important. In [42], attention maps are produced to focus computation on specific parts of the input. The weighted attention pooling layer is plugged in as a replacement for a pooling operation in a fully convolutional network. In [202], a three-stream architecture is proposed which includes two attention streams and a global pooling stream. A shared ResNet is used to extract spatial features for all three streams. Each attention layer employs a fusion layer to combine global and local information and produces composite features. Furthermore, global-attention regularization is proposed to guide two attention streams to better model dynamics of composite features with the reference to the global information.

**Sparse Features.** Sparse features, particularly skeletons, are very popular in action recognition due to their action-focusing nature and compactness. Several studies use human skeleton information as a sequence of joint coordinate lists [79, 187, 196] where the coordinates are extracted by pose estimators.

While using CNN networks to process RGB frames is straightforward, in skeleton-based action modeling, network is faced with the challenge of arranging skeleton features. Earlier methods [79, 196] simply use the keypoint coordinates to generate a sequence of feature vectors. The issue with this method is that it focuses on modeling

the information in the time domain and doesn't explore the spatial relations between body joints. Other approaches arranged skeleton data as a pseudo-image [68, 70, 77]. However, recent works have shown that graph networks can efficiently model non-Euclidean data like human skeletons. Performance results from table 8 shows the superior performance of arranging skeleton data in a graph structure.

**Table 8** Comparison of skeleton-based action recognition performance for NTU RGB-D [131] dataset with different feature arrangements. mAP refers to mean average precision.

model	[32]	[131]	[99]	[70]	[68]	[187]	[169]	[80]
arrangement	vector	vector	vector	pseudo-image	pseudo-image	graph	<b>graph</b>	<b>graph</b>
mAP	59.1%	62.9%	69.2%	74.3%	79.6%	81.5%	<b>84.2%</b>	<b>86.8%</b>

With introduction of ST-GCN in [187], spatio-temporal graph convolutions became a research hotspot. The core of this approach is to integrate temporal module in the spatial GCN. Several variants of ST-GCN are proposed [20, 21, 170] to improve the network capacity and computation consumption of the network. However, the main limitation of ST-GCN is that it ignores the semantic connections among intra-frame joints by using a fixed graph structure. Recognizing action lies in looking beyond the local joint connectivity as learning not only happens in spatially connected joints, but also in the potential dependence of disconnected joints. For example, in “walking” there is a high correlation between arms and legs while they are spatially apart. This achieves by extracting multi-scale features and long-range temporal dependencies, as joints that are spatially apart can also have strong correlations [106]. In this regard, some techniques have been adopted to enhance the flexibility of GCNs. Attempts from using adaptive learning graph structure in [133, 136] to designing a graph-based search space to explore spatio-temporal connections [122] has been made. In [27], dilated convolutions [192] are adopted to increase receptive field size and capture multi-scale context without increasing model complexity. In [106] this issue is addressed by performing graph convolutions with higher-order polynomials of the skeleton adjacency matrix which increases the receptive field of graph convolutions. Attention mechanisms are also adopted to improve the ability of extracting high-level joints. In [136] a spatio-temporal channel attention module is embedded in each layer of the GCN, which enables model to focus on the discriminative details of joints, frames and channels in action recognition.

*Pros.* When sparse features are used, since only pose information is included, they contain high-level semantic information in a small amount of data and are more robust in dynamic circumstances [59, 135]. As skeleton data does not contain color information, it is not affected by the limitations of RGB frames [138], and can provide a stable low-frequency representation of human actions.

*Cons.* There are some challenges in modeling sparse feature representation. A first limitation is the “reliability” which means the recognition ability of sparse features is mainly affected by the distribution shift of coordinates. As joint coordinates are produced by a pose estimator, applying a different pose estimation algorithm may lead to a small perturbation of coordinates which causes different predictions [206].

Also, local subtle motion patterns are lost in the process of pose estimation. Therefore, sparse nature of skeleton sequences is sometimes not informative enough for describing subtle actions like “reading”, “writing”, and “shaking head”.

Another challenge is the lack of scalability of sparse-features. As sparse features are defined for every human separately. For example, each joint of human skeleton is defined as a node per person, the complexity of network linearly increases with increasing the number of persons, which limits its applicability in multi-person or group activity recognition.

*Solutions.* To overcome reliability issues, many approaches take advantage of multi-modal visual resources including RGB frames, depth maps and joint heatmaps to compensate for the lack of information in local and global domain [11, 33]. People use other forms of representations using heatmap volumes to show skeletons to alleviate the scalability issue of sparse features [24, 185].

**Frame-level Features.** Some approaches extract the motion cues between adjacent frames and learn frame-level temporal dependencies in a sequential signature. In action recognition, modeling both short-range and long-range motions is sometimes required. (1) To this end, some earlier methods firstly extract hand-crafted optical flow [139, 161], then use a 2D CNN-based two-stream framework to process optical flow and RGB frames in each stream separately. These lines of works have several drawbacks: First, computing optical flow is time-consuming and storage demanding. Furthermore, the training of spatial and temporal features is separated, and the fusion of two streams is performed only at the late layers. Several following works had improved this framework by using different mid-level links to fuse the features of two separated streams [37, 38]. However, these methods still require additional time and storage costs for computation of optical flow. (2) Another line of work aggregates temporal information by sequence learning [31, 89]. The majority of these methods treat each frame, or point in time, with equal weight, but not all parts of the video are equally important and thus it is also key that we develop feature extraction methods that can determine where to extract features from. First, non-uniformly extracting features efficiently from only informative temporal episodes is challenging as it is required to look at the whole video to determine which parts are informative. Some recent work [74, 113, 174] have proved that a recognition system can benefit from selecting the informative frames rather than simply taking the uniformly sampled frames as inputs. However, these systems treat the frame selection and feature extraction as two separate stages and thus the frame selection can not benefit from the later feature extraction thus reducing the descriptive power of the network and adding redundancy in the two stages. In order to tackle this challenge, in [84], a two-branch architecture is suggested that maintains both uniformly frame-level features and non-uniformly chunk-level features in an end to end manner. To produce non-uniformly features, a temporal map is used that non-uniformly projects temporal instances to a smaller subset by using self-attention-like module. This component is proposed to only sample the most informative frames across time. (3) Another recent line of work adopts video transformers that apply self-attention to spatial-temporal features. Representative networks include TimeSformer [6], ViViT [1], Mformer [121] and MViT [36]. Combining 2D backbones and Transformers, VTN [115] and CARL [14] can efficiently

process long video sequences. However, these networks are designed to process a batch of frames at once on video tasks which requires large computing memory. In [190], a recursive mechanism is deployed to process the videos frame by frame and consume less GPU memory.

*Pros.* Generally, standard architectures of frame-level features have lower computational costs compared to chunk-level features. Also, frame-level features are more concise when aggregating local frame-level features are required for a global compact representations.

*Cons.* The main challenge in extracting frame-level features is “co-occurrence representation” learning which in action recognition refers to when a model needs to learn a set of human actions with a specific set of spatial features at certain times. For example, in the action of “touching back”, model needs to focus first on the hand and then pay attention to the back [83]. In many of existing approaches, a temporal module and a spatial module are designed separately. Their typical approach is to use a convolutional network to extract spatial relations in each frame, then use a 1D convolution [80, 133, 134] or LSTM [76, 82] to model temporal dependencies. However, such decoupled design restricts the direct information flow across space-time for capturing complex regional spatial-temporal joint dependencies.

*Solutions.* To help model better learn co-occurrence representation, message passing and cross connection strategy is adopted to avoid stacking multiple spatio-temporal modules and transfer information. In [189], the feedback connection was integrated into GCN to transfer the high-level semantic features to the low-level layer, and gradually transmitted the temporal information to build the global spatio-temporal action recognition model. In [150], attentions provided from temporal-stream is used to help spatial stream by cross-link layers. In [83], a coordinate system conversion and spatio-temporal-unit feature enhancement is proposed to perform co-occurrence learning. In [138], each joint coordinate is transformed into a spatial feature with a linear layer. Then data is augmented with the spatial feature and the difference between spatial features between consecutive frames. Then a shared LSTM and three layers of graph convolutional LSTM are applied to model co-occurrence representation learning between joints.

**Chunk-level Features.** Another type of approaches is to extract chunk-level temporal features in a global signature. In this case, the chunk is defined as multi-dimensional time series of dense/sparse features. Thanks to the ability of 3D CNNs to implicitly model motion information along with the semantic features, this line of works has seen significant advances in recent years. The first work in this line was C3D [151], which proposed using 3D convolutions to jointly model the spatial and temporal features in a global signature. To use pre-trained 2D CNNs, in [12], I3D was proposed that inflates the pre-trained 2D convolutions to 3D ones.

*Pros.* Generally, chunk-level features benefit from the co-occurrence representation learning as there is a link between temporal and spatial channels. These networks are potentially more effective in learning fine detailed actions such as “sitting” and “standing up”.

*Cons.* While so many attempts have been done in capturing motion using chunk-level features, most of approaches often lack specific consideration in the temporal

dimension. Therefore, designing an effective temporal module of high motion modeling power and low computational consumption is still a challenging problem. First, the 3D networks require a substantial amount of computation and time. Also, compared to 2D kernels, 3D convolutions have to reduce the spatial resolution to decrease memory consumption, which may lead to the loss of finer details.

Moreover, temporal features are typically extracted from a fixed-length clip instead of a length-adaptive clip which is not suited for different visual tempos. Visual tempo defines the speed of an action, meaning how fast and slow an action is performed at the temporal scale which in some cases is crucial for recognizing actions e.x. walking, jogging and running.

*Solutions.* To decrease the heavy computations of 3D CNNs, some works proposed to factorize the 3D convolution with a 2D spatial convolution and a 1D temporal convolution [45, 94, 152] or a mixed up of 2D CNN and 3D CNN [176, 208]. In [152], a group convolution is used to disentangle channel interactions and spatio-temporal interactions, or use separated channel groups to encode both spatial and spatio-temporal interactions in parallel with 2D and 3D convolution [109]. While all these existing approaches are designed to deal with the curse of dimension, there is a lack of data dependent decision to adaptively guide features through different routs in the network. In [162], two temporal difference module is proposed which computes multi-scale and bidirectional motion information between frames and chunks. In [146], features are selectively routed through temporal dimension and are combined with each other without any computational overhead.

Some people use multi-scale convolutional kernels to cover various visual tempos. In [188], multi-scale convolutional features are incorporated in asymmetric 3D convolutions to improve temporal feature learning capacity. In [58], Timeception layer is designed which temporally convolves each chunk using multi-scale temporal convolution module to tolerate a variety of temporal extents in a complex action. Some people design a level-specific network frame pyramid to handle the variance of visual tempos [39, 193]. In [26], a multi-scale transformer is proposed which is built on top of temporal segments using 3D convolutions in a token-based architecture to promote multiple temporal scales of tokens. Having different scales allows model to capture both fine-grained and composite actions across time. In [153], a direct attention mechanism is incorporated in transformers to exploit the direction of attention across frames and correct the incorrectly-ordered frames to the right ones and provide an accurate prediction.

**Summary.** We summarized pros and cons of different features and the possible solution in table 7. Dense features have the advantage of using appearance info in action recognition while they suffer from background noise and representation bias. Possible solutions for background noise and representation bias include augmenting additional information to the input and well calibrated dataset and attention mechanism, respectively. Sparse features are more robust against background noise and have lower computational cost compared to dense features, while they suffer from lack of reliability and scalability. Possible solutions for their drawbacks could be augmenting additional information and visual sources. Researches have shown in spatial domain, using multi-modal inputs, particularly accompanied with attention mechanism are

**Table 9** Comparison of different action recognition performance for NTU RGB-D [131] dataset. CS and CV refer to cross subject (various human subjects) and cross view (various camera views) split of the dataset.

Model	Spatial	Temporal	CS accuracy	CV accuracy
[4]	Dense	Frame	86.6%	93.2%
[207]	Dense	Chunk	94.3%	97.2%
[10]	Sparse	Chunk	76.5%	84.7%
[148]	Sparse	Chunk	87.5%	93.2%
[28]	Dense + Sparse	Chunk	91.8%	94.9%
[103]	Dense + Sparse	Frame	91.7%	95.2%
[145]	Dense + Sparse	Frame	92.2%	-
[8]	Dense + Sparse	Chunk	92.5%	97.4%
[33]	Dense + Sparse	Chunk	97.0%	99.6%
[9]	Dense + Sparse	Chunk	96.0%	98.8%

very helpful in understanding human actions, as humans’ brain adopt all visual inputs to recognize an action. The quantitative results are shown in table 9 on NTU RGB-D dataset confirm that using multi-modal spatial features outperforms single-modal approaches.

In temporal domain, using frame-level features has the advantage of lower computational cost, compared to chunk-level features. However, decoupling spatial and temporal dimensions restricts co-occurrence representation learning in distinguishing complex actions. Some studies alleviate this problem by message passing techniques and cross-link connections. On the contrary, chunk-level features allow establish of connections and links between temporal and spatial dimensions to learn order and co-occurrence of micro-actions. However, typically these models take fixed-length instead of adaptive-length chunks. Some works address this issue by using frame pyramids and multi-scale convolutions. Another drawback of chunk-level features is their high computational costs which could be alleviated by disentangling convolutions in different layers of network. To conclude, as shown in table 7, in presence of sparse features, using chunk-level approaches outperforms frame-wise methods due to reducing number of parameters, increasing depth of network and co-occurrence representation learning.

### 3.2 Video Object Segmentation

Video object segmentation (VOS) is a video processing technique. The goal of VOS is to segment pixel-level masks of foreground objects in every frame of a given video. VOS has attracted extensive attention these years because it can be applied to diverse fields in computer vision. Recent VOS research can be divided into two sub-tasks: semi-supervised and unsupervised. The semi-supervised VOS aims to re-locate and segment one or more objects that are given in the first frame of a video in pixel-level masks. The unsupervised VOS aims to automatically segment the object of interest from the background, usually, the most salient object(s) will be segmented. Generally, the input to the video object segmentation is a sequence of color frames in RGB format. Feature modeling is the first stage of these video object segmentation pipelines. In the following, we discuss pros and cons of different features in VOS.

**Dense Features.** Dense features are widely used because of the neutrality of VOS which is supposed to estimate the pixel-level object masks. These methods often apply

**Table 10** Pros and Cons of spatial and temporal features in the video object segmentation (VOS) task. There are some abbreviations in the table: ap. and info. means appearance and information, respectively.

Features	Pros	Cons	Solution	References
Dense	contains rich ap. info.	weak in occlusion-handling.	occlusion-aware encoders	[66, 67]
Sparse	fast & low comp.	less accuracy	hybrid algorithm	[16, 147]
Frame	low computation cost low latency	error accumulation past frame management	dynamic feature space management	[81, 92]
Chunk	multi-modal modeling	large computation cost	knowledge distillation	[171, 172]

a pre-trained CNN encoder to extract dense feature maps from each frame. Generally, 2D CNN encoder is widely-used to extract feature maps [2, 96, 98, 120]. The CNN encoder is often pretrained on ImageNet [30] and fine-tuned on the video object segmentation dataset. After extracting the dense feature maps, these methods utilize a transformer [155] to encode the feature maps into two *keys* and *values*, where *keys* contains the semantic code of the object and *values* contains the detailed appearance information. The encoded feature maps can be used for matching and information retrieval. Besides the appearance feature maps, Zhang et al. [199] introduced perceptual consistency to aid with predicting the pixel-wise correctness of the segmentation on an unlabeled frame.

*Pros.* RGB frame provides rich information about the appearances textcolorgreen-which is crucial in VOS. With the help of GPU parallel computing and large-scale pretrained model, dense features extraction from RGB frame becomes standard processing in video object segmentation nowadays.

*Cons.* Although dense feature extraction has been widely studied, They may contain background noise and irrelevant information. Particularly, occlusion-handling is a main challenge of using dense features in VOS.

*Solutions.* In dealing with occlusions, its imperative to recover occluded parts by corresponding shape and appearance in motion rather than irrelevant background. For this purpose, Occlusion-aware feature modeling [67] is proposed which uses shape completion and flow completion modules to inpaint invisible parts intelligently. In [66], a pipeline is proposed that used GCN which allows propagation of non-local information across pixels despite the presence of occluding regions.

**Sparse Features.** Compared with the dense feature extraction, sparse feature modeling discards the irrelevant information from input. In VOS, short tracks or tracklets in a frame are considered as sparse features. In [16], a State-Aware Tracker (SAT) is proposed that takes advantage of the inter-frame consistency and deal with each target object as a tracklet. Because the irrelevant background is discard, they achieve real-time video object segmentation on 39 FPS which is faster than traditional dense feature methods. In [147], the video object tracking module is adopted to first locate the object region from the background. Then they segment the object masks from the small object region.

*Pros.* The advantages of the sparse feature modeling are two fold. Firstly, it uses pre-processing or prior knowledge to clean the irrelevant information from input, which

makes the pipeline more robust to the noises from the background. Secondly, using sparse features can accelerate the speed of segmentation, which gives the users an option to choose a trade-off between efficiency and accuracy.

*Cons.* Sparse feature modeling in VOS relies on the data pre-processing. It is inevitable that some important data from the object is mistakenly discard in this process. Also, sparse features lack the appearance/color information which is helpful in VOS.

*Solution.* Compared with dense feature modeling, sparse feature modeling has lower accuracy but better runtime performance. Recently, [173] used tracklet query and tracklet proposal that combines RoI features and dense frame features by the vision transformer. It achieved state-of-the-art performance in YouTube-VIS 2019 val set [191].

**Table 11** Comparison of different spatial feature performance for DAVIS-2016 and DAVIS-2017 datasets [125].  $\mathcal{J}$ ,  $\mathcal{F}$  and  $\mathcal{J}\&\mathcal{F}$  score represent region similarity, contour accuracy and the average value of region similarity and contour accuracy, respectively. All models use frame-level features in temporal dimension.

Model	Spatial	DAVIS val 16			DAVIS val 17			FPS
		$\mathcal{J}$ (%)	$\mathcal{F}$ (%)	$\mathcal{J}\&\mathcal{F}$ (%)	$\mathcal{J}$ (%)	$\mathcal{F}$ (%)	$\mathcal{J}\&\mathcal{F}$ (%)	
FEELVOS [158]	Dense	81.1	82.2	81.7	69.1	74.0	71.5	2.2
STM [117]	Dense	84.8	88.1	86.5	69.2	74.0	71.6	6.3
AGAME [62]	Dense	82.0	-	-	67.2	72.7	70.0	14.3
AGSS-VOS [93]	Dense	-	-	-	63.4	69.8	66.6	10.0
AFB-URR [92]	Dense	-	-	-	73.0	76.1	74.6	4.0
FRTM [129]	Dense	-	-	81.7	66.4	71.2	68.8	21.9
LCM [53]	Dense	-	-	-	73.1	77.2	75.2	8.5
RMNet [175]	Dense	80.6	82.3	81.5	72.8	77.2	75.0	11.9
SWEM [96]	Dense	<b>87.3</b>	<b>89.0</b>	<b>88.1</b>	<b>74.5</b>	<b>79.8</b>	<b>77.2</b>	36
BMVOS [22]	Dense	82.9	81.4	82.2	70.7	74.7	72.7	45.9
FTM [147]	Sparse	77.5	-	78.9	69.1	-	70.6	11.1
SAT [16]	Sparse	82.6	83.6	83.1	68.6	76.0	72.3	39.0
SAT-Fast [16]	Sparse	-	-	-	65.4	73.6	69.5	<b>60.0</b>

**Frame-level Features.** Several VOS methods predict segmentation masks frame by frame. They utilize previous frames information to model the current object appearance and motion. Appearance similarity between frames is widely-used to segment objects from background [2, 62, 81, 116, 168, 180, 199, 205]. The intuition behind these papers is that perceptually similar pixels are more likely to be in the same class. In frame-level features, it is essential to make full use of historical frames in the videos. Some methods use space-time memory banks to store the embeddings every several frames [19, 53]. Other methods fuse pixel embedding of the current frame and the memory bank [92].

In [157] a recurrent pipeline (RVOS) is proposed to keep the coherence of the segmented objects along time. Due to the RNN’s memory capabilities, RVOS is recurrent in the spatial-temporal domain and can handle the instances matching at different frames. Recent papers used auxiliary temporal information including optical flow [3, 93, 195] to aid VOS.

**Table 12** Comparison of different spatial feature performances for Youtube-VOS [182] dataset.  $\mathcal{J}$  and  $\mathcal{F}$  represent region similarity, contour accuracy and seen and unseen tags indicate seen and unseen objects from the training dataset separately.  $\mathcal{G}$  measures the overall score of the segmentation accuracy. FPS metric reports the runtime speed in frames per second.

Model	Spatial	$\mathcal{J}_{seen}(\%)$	$\mathcal{J}_{unseen}(\%)$	$\mathcal{F}_{seen}(\%)$	$\mathcal{F}_{unseen}(\%)$	$\mathcal{G}$	FPS
STM [117]	Dense	79.7	84.2	72.8	80.9	79.4	-
RVOS [157]	Dense	63.6	45.5	67.2	51.0	56.8	22.7
AGSS-VOS [93]	Dense	71.3	65.5	76.2	73.1	71.3	12.5
AFB-URR [92]	Dense	78.8	74.1	<b>83.1</b>	82.6	79.6	3.3
FRTM [129]	Dense	72.3	65.9	76.2	74.1	72.1	-
LCM [53]	Dense	82.2	75.7	86.7	83.4	82.0	-
RMNet [175]	Dense	82.1	85.7	75.7	82.4	81.5	-
SWEM [96]	Dense	<b>82.4</b>	<b>86.9</b>	77.1	<b>85.0</b>	<b>82.8</b>	-
BMVOS [22]	Dense	73.5	68.5	77.4	76.0	73.9	28.0
SAT [16]	Sparse	67.1	55.3	70.2	61.7	63.6	<b>39.0</b>

*Pros.* Most VOS methods follow the frame-level feature modeling because it usually has lower computational costs. These features are also suitable to handle streaming data such as online video or online meeting because it processes the video frame by frames. Using optical flow map as dense map is very common in the unsupervised video object segmentation task [128, 194], which provide the dense correspondence between similar pixels.

*Cons.* Although frame-level feature modeling has lots of benefits, it brings some shortcomings and limitations. First of all, as information is processed frame by frame, past frames aid model in segmenting objects in the current and future frames. Managing past information is very challenging and requires lots of efforts in model design, particularly when the length of video increases. Additionally, the errors from past frames are accumulated and transferred during video processing.

*Solutions.* Recent research papers proposed adaptive memory bank scheme to maintain features from previous frames. In [92] an exponential moving average (EMA) based scheme is proposed to update the record of historical information. This method helps model to better use past information. In [81], a learning based spatial-temporal aggregation model (SAM) is introduced to distill the frame-level features and automatically correct the accumulated errors.

**Chunk-level Features.** Inspired by the human action recognition field, In [55], a 3D convolutional neural network is used as a backbone to extract chunk-level features. Recent chunk-level feature modeling is used in the multimodal video segmentation task, which combines text reasoning, video understanding, instance segmentation and tracking. Recent papers first use a CNN encoder to extract feature maps from each frame [127], then combine them into the spatial-temporal features. In [7, 171, 172], a Transformer is used to model spatial-temporal features from videos.

*Pros.* The advantage of chunk-level features in VOS is that it models the global semantic information in consecutive frames, which is critical to multi-modal tasks.

*Cons.* Modeling chunk-level features often requires large memory and computation costs compared with frame-level features. Heavy computation requirement limits its application on mobile device.

*Solutions.* Knowledge distillation [95] and neural network pruning/search [204] are two active fields aim to relieve the burden of neural network computing. They can help in finding a smaller network with lower network parameters but similar accuracy performance. These methods could be used to optimize the architecture of the chunk-level feature encoder to accelerate their computation. In [179], an acceleration framework is proposed based on video-compressed codec. For each chunk, it has three types of I-frames, P-frames, and B-frames. They utilize the sparsity of I-frames and motion vectors from P-/B-frames to accelerate the video object segmentation.

**Table 13** Comparison of different models performance with chunk level features for Ref-Youtube-VOS [130] across different backbones and chunk sizes  $\omega$ .  $\mathcal{J}$ ,  $\mathcal{F}$  and  $\mathcal{J}\&\mathcal{F}$  score represent region similarity, contour accuracy and the average value of region similarity and contour accuracy, respectively.

Model	Backbone	$\mathcal{J}$ (%)	$\mathcal{F}$ (%)	$\mathcal{J}\&\mathcal{F}$ (%)
URVOS [130]	ResNet-50	45.3	49.2	47.2
MLRL [171]	ResNet-50	48.4	51.0	49.7
MTTR ( $\omega=12$ ) [7]	Video-Swin-T	54.0	56.6	55.3
ReferFormer ( $\omega=5$ ) [172]	Video-Swin-T	54.8	57.3	56.0
ReferFormer [172]	Video-Swin-B	<b>61.3</b>	<b>64.6</b>	<b>62.9</b>

**Table 14** Comparison of different spatial feature performance for A2D-Sentences [41].  $*^T$  indicates the backbone architecture is Video-Swin-T.  $*^B$  indicates the backbone architecture is Video-Swin-B. The Precision @K measures the percentage of test samples that their whole IoU scores are higher than threshold K.

Model	Precision					IoU		mAP
	P@0.5	P@0.6	P@0.7	P@0.8	P@0.9	Overall	Mean	
MTTR <sup>T</sup> ( $\omega=8$ ) [7]	72.1	68.4	60.7	45.6	16.4	70.2	61.8	44.7
MTTR <sup>T</sup> ( $\omega=10$ ) [7]	75.4	71.2	63.8	48.5	16.9	72.0	64.0	46.1
ReferFormer <sup>T</sup> ( $\omega=6$ ) [172]	76.0	72.2	65.4	49.8	17.9	72.3	64.1	48.6
ReferFormer <sup>B</sup> ( $\omega=5$ ) [172]	<b>83.1</b>	<b>80.4</b>	<b>74.1</b>	<b>57.9</b>	<b>21.2</b>	<b>78.6</b>	<b>70.3</b>	<b>55.0</b>

**Summary.** We summarized pros and cons of features and the possible solution in table 10. Dense features Contain rich appearance information which is essential in VOS, but they are not robust against occlusion. Some approaches use occlusion-aware encoders to overcome this shortcoming. Sparse features disregard irrelevant information and have lower computation cost. The accuracy of sparse features are lower in VOS task. To overcome this issue, people use hybrid algorithms. Frame-level features have lower computational cost compared to frame-level features and are suitable for stream processing. However, the computation error is accumulated in these features while using past information. Also, managing past frames requires a lot of efforts. Some researches adopted dynamic memory to relieve the error accumulation. Also they used exponential weight smoothing and learning based feature adaption to manage past frame information. Chunk-level features model more global information

than frame-level features but are more computationally expensive. Different methods are adopted that use knowledge distillation to alleviate computation burden.

To conclude as shown in Table 11 and Table 12, we compared the performance of different spatial feature models. These comparisons are from the spatial aspect, they all belong to the frame-level features. The DAVIS’17 datasets [125] benchmarks contain 60 videos for training and 30 videos for validation. The YouTube-VOS [182] benchmark contains 3,471 videos for training and 507 for validation. In both benchmarks, the video object segmentation task is to segment and track the arbitrary number of objects in each video. The groundtruth masks for each object are provided in the first frame. For DAVIS datasets, we followed the official evaluation metrics, region similarity  $\mathcal{J}$  and contour accuracy  $\mathcal{F}$ . The  $\mathcal{J}\&\mathcal{F}$  score is the average value of region similarity and contour accuracy. For YouTube-VOS dataset, we used similar metrics  $\mathcal{J}$  and  $\mathcal{F}$ . *seen* and *unseen* tags indicate seen and unseen objects from the training dataset separately.  $\mathcal{G}$  measures the overall score of the segmentation accuracy. FPS metric reports the runtime speed in frames per second. The runtime speed was measured on NVIDIA 2080Ti and NVIDIA V100 GPUs. SWEM [96] has the best accuracy performance in both DAVIS and YouTube-VOS datasets because it utilizes dense feature modeling, while SAT [16] achieves the best runtime speed because of sparse feature modeling.

To compare the chunk-level feature modeling, we chose the popular Ref-YouTube-VOS [130] and A2D-Sentences [41] benchmarks for comparisons. The Ref-YouTube-VOS dataset covers 3,978 videos with around 15K language descriptions. We used similar metrics to measure the segmentation performance, region similarity  $\mathcal{J}$ ,  $\mathcal{F}$ , and  $\mathcal{J}\&\mathcal{F}$ . The A2D-Sentences dataset contains 3,782 videos and each video has 3-5 frames annotated with the pixel-level segmentation masks. The model is evaluated with criteria of Precision @K, Overall IoU, Mean IoU and mAP. The Precision @K measures the percentage of test samples whose IoU scores are higher than threshold K. Following standard protocol, the thresholds are set as 0.5:0.1:0.9. We compared the performance results across different network settings: (1) different backbones Video-Swin-T and Video-Swin-B from Video Swin Transformer [107], (2) different chunk sizes  $\omega$ . From Table 13 and Table 14, ReferFormer [172] achieves the best performance because it designed the cross-model feature pyramid network to extract multi-scale chunk-level features from the input.

## 4 Conclusion and Future Work

We provided an extensive survey on recent studies on deep video representation learning. We provided a new taxonomy of these features and classify existing methods accordingly. We discussed and compared the effectiveness (robustness) of different types of features under scenarios with different types of noise.

**Challenges.** Spatially dense features can encode rich contextual information but are more sensitive to background noise. Handling dense features in presence of intense occlusion and view variations is still challenging. In contrast, sparse features are more robust against background noise and illumination variance, but arranging sparse topologies in spatial dimension is still challenging. It is still an open question if the spatial relations should be defined based on the natural inherent relation of the features

points or on the correlation of non-adjacent points throughout the video. For example, in human body skeleton, should the spatial relations between coordinates be defined based on natural relations of human body joints or the correlation of non-adjacent joints during the movement? In frame-level features, lack of cross connections between temporal and spatial domains is a major drawback in capturing complex dynamics. In chunk-level features, improving model’s generalizability and high computational cost are the main challenges.

**Future directions.** The drawbacks of either using sparse or dense features could be solved by using multi modal inputs to some extent. A future direction of these studies is on designing new methods for mapping between different modalities’ feature space, learning effective representations from multiple data modalities, and understanding when and where the fusion should happen. Recent studies [209] demonstrated using proper multi-modals clearly improve video analysis performance. While current attention methods have achieved progress in video representation learning, they often bring higher model complexity and suffer from heavier computational burden. Hence, many recent studies are on building more efficient attention models [50].

#### *Data availability statement*

All data supporting the findings of this study are available within the paper.

## Declarations

**Competing interests** We do not have any conflict of interest related to the manuscript.

## References

- [1] Arnab A, Dehghani M, Heigold G, et al (2021) Vivit: A video vision transformer. In: ICCV, pp 6836–6846
- [2] Athar A, Luiten J, Hermans A, et al (2022) Hodor: High-level object descriptors for object re-segmentation in video learned from static images. In: CVPR, pp 3022–3031
- [3] Azulay A, Halperin T, Vantzos O, et al (2022) Temporally stable video segmentation without video annotations. In: WACV, pp 3449–3458
- [4] Baradel F, Wolf C, Mille J, et al (2018) Glimpse clouds: Human activity recognition from unstructured feature points. In: CVPR, pp 469–478
- [5] Bendre N, Zand N, Bhattarai S, et al (2022) Natural disaster analytics using high resolution satellite images. In: World Automation Congress, IEEE, pp 371–378
- [6] Bertasius G, Wang H, Torresani L (2021) Is space-time attention all you need for video understanding? In: ICML, p 4

- [7] Botach A, Zheltonozhskii E, Baskin C (2022) End-to-end referring video object segmentation with multimodal transformers. In: CVPR, pp 4985–4995
- [8] Bruce X, Liu Y, Chan KC (2021) Multimodal fusion via teacher-student network for indoor action recognition. In: AAAI, pp 3199–3207
- [9] Bruce X, Liu Y, Zhang X, et al (2022) Mmnet: A model-based multimodal network for human action recognition in rgb-d videos. PAMI
- [10] Caetano C, Sena J, Brémond F, et al (2019) Skelemotion: A new representation of skeleton joint sequences based on motion information for 3d action recognition. In: International conference on advanced video and signal based surveillance, IEEE, pp 1–8
- [11] Cai J, Jiang N, Han X, et al (2021) Jolo-gcn: mining joint-centered light-weight information for skeleton-based action recognition. In: WACV, pp 2735–2744
- [12] Carreira J, Zisserman A (2017) Quo vadis, action recognition? a new model and the kinetics dataset. In: CVPR, pp 6299–6308
- [13] Chen D, Li H, Xiao T, et al (2018) Video person re-identification with competitive snippet-similarity aggregation and co-attentive snippet embedding. In: CVPR, pp 1169–1178
- [14] Chen M, Wei F, Li C, et al (2022) Frame-wise action representations for long videos via sequence contrastive learning. In: CVPR, pp 13801–13810
- [15] Chen X, Yuille AL (2015) Parsing occluded people by flexible compositions. In: CVPR, pp 3945–3954
- [16] Chen X, Li Z, Yuan Y, et al (2020) State-aware tracker for real-time video object segmentation. In: CVPR, pp 9384–9393
- [17] Chen Y, Pont-Tuset J, Montes A, et al (2018) Blazingly fast video object segmentation with pixel-wise metric learning. In: CVPR, pp 1189–1198
- [18] Chen Z, Wang X, Sun Z, et al (2016) Motion saliency detection using a temporal fourier transform. *Optics & Laser Technology* 80:1–15
- [19] Cheng HK, Tai YW, Tang CK (2021) Modular interactive video object segmentation: Interaction-to-mask, propagation and difference-aware fusion. In: CVPR, pp 5559–5568
- [20] Cheng K, Zhang Y, Cao C, et al (2020) Decoupling gcn with dropgraph module for skeleton-based action recognition. In: ECCV, Springer, pp 536–553
- [21] Cheng K, Zhang Y, He X, et al (2020) Skeleton-based action recognition with shift graph convolutional network. In: CVPR, pp 183–192

- [22] Cho S, Lee H, Kim M, et al (2022) Pixel-level bijective matching for video object segmentation. In: WACV, pp 129–138
- [23] Choi J, Gao C, Messou JC, et al (2019) Why can't i dance in the mall? learning to mitigate scene bias in action recognition. NIPS 32
- [24] Choutas V, Weinzaepfel P, Revaud J, et al (2018) Potion: Pose motion representation for action recognition. In: CVPR, pp 7024–7033
- [25] Cuevas C, Quilón D, García N (2020) Techniques and applications for soccer video analysis: A survey. *Multimedia Tools and Applications* 79(39-40):29685–29721
- [26] Dai R, Das S, Kahatapitiya K, et al (2022) Ms-tct: Multi-scale temporal convtransformer for action detection. In: CVPR, pp 20041–20051
- [27] Dai X, Singh B, Ng JYH, et al (2019) Tan: Temporal aggregation network for dense multi-label action recognition. In: WACV, IEEE, pp 151–160
- [28] De Boissiere AM, Noumeir R (2020) Infrared and 3d skeleton feature fusion for rgb-d action recognition. *IEEE Access* 8:168297–168308
- [29] Deng J, Dong W, Socher R, et al (2009) Imagenet: A large-scale hierarchical image database. In: CVPR, pp 248–255, <https://doi.org/10.1109/CVPR.2009.5206848>
- [30] Deng J, Dong W, Socher R, et al (2009) Imagenet: A large-scale hierarchical image database. In: CVPR, Ieee, pp 248–255
- [31] Donahue J, Anne Hendricks L, Guadarrama S, et al (2015) Long-term recurrent convolutional networks for visual recognition and description. In: CVPR, pp 2625–2634
- [32] Du Y, Wang W, Wang L (2015) Hierarchical recurrent neural network for skeleton based action recognition. In: CVPR, pp 1110–1118
- [33] Duan H, Zhao Y, Chen K, et al (2022) Revisiting skeleton-based action recognition. In: CVPR, pp 2969–2978
- [34] Eun H, Moon J, Park J, et al (2020) Learning to discriminate information for online action detection. In: CVPR, pp 809–818
- [35] Fabbri M, Lanzi F, Calderara S, et al (2018) Learning to detect and track visible and occluded body joints in a virtual world. In: ECCV
- [36] Fan H, Xiong B, Mangalam K, et al (2021) Multiscale vision transformers. In: ICCV, pp 6824–6835

- [37] Feichtenhofer C, Pinz A, Zisserman A (2016) Convolutional two-stream network fusion for video action recognition. In: CVPR, pp 1933–1941
- [38] Feichtenhofer C, Pinz A, Wildes RP (2017) Spatiotemporal multiplier networks for video action recognition. In: CVPR, pp 4768–4777
- [39] Feichtenhofer C, Fan H, Malik J, et al (2019) Slowfast networks for video recognition. In: ICCV, pp 6202–6211
- [40] Gao R, Oh TH, Grauman K, et al (2020) Listen to look: Action recognition by previewing audio. In: CVPR, pp 10457–10467
- [41] Gavriluk K, Ghodrati A, Li Z, et al (2018) Actor and action video segmentation from a sentence. In: CVPR, pp 5958–5966
- [42] Girdhar R, Ramanan D (2017) Attentional pooling for action recognition. *Advances in Neural Information Processing Systems* 30
- [43] Hamilton WL, Ying R, Leskovec J (2017) Representation learning on graphs: Methods and applications. arXiv preprint arXiv:170905584
- [44] Hao X, Li J, Guo Y, et al (2021) Hypergraph neural network for skeleton-based action recognition. *TIP* 30:2263–2275
- [45] He D, Zhou Z, Gan C, et al (2019) Stnet: Local and global spatial-temporal modeling for action recognition. In: AAAI, pp 8401–8408
- [46] He K, Zhang X, Ren S, et al (2016) Deep residual learning for image recognition. In: CVPR, pp 770–778
- [47] He K, Gkioxari G, Dollár P, et al (2017) Mask r-cnn. In: ICCV, pp 2961–2969
- [48] Herzig R, Ben-Avraham E, Mangalam K, et al (2022) Object-region video transformers. In: CVPR, pp 3148–3159
- [49] Horn BK, Schunck BG (1981) Determining optical flow. *Artificial intelligence* 17(1-3):185–203
- [50] Hou Q, Zhou D, Feng J (2021) Coordinate attention for efficient mobile network design. In: CVPR, pp 13713–13722
- [51] Hou R, Ma B, Chang H, et al (2019) Vrstc: Occlusion-free video person re-identification. In: CVPR, pp 7183–7192
- [52] Hu JF, Zheng WS, Lai J, et al (2015) Jointly learning heterogeneous features for rgb-d activity recognition. In: CVPR, pp 5344–5352

- [53] Hu L, Zhang P, Zhang B, et al (2021) Learning position and target consistency for memory-based video object segmentation. In: CVPR, pp 4144–4154
- [54] Hu YT, Huang JB, Schwing AG (2018) Videomatch: Matching based video object segmentation. In: ECCV, pp 54–70
- [55] Huang X, Xu J, Tai YW, et al (2020) Fast video object segmentation with temporal aggregation network and dynamic template matching. In: CVPR, pp 8879–8889
- [56] Huang Z, Wan C, Probst T, et al (2017) Deep learning on lie groups for skeleton-based action recognition. In: CVPR, pp 6099–6108
- [57] Hussain T, Muhammad K, Ding W, et al (2021) A comprehensive survey of multi-view video summarization. *Pattern Recognition* 109:107567
- [58] Hussein N, Gavves E, Smeulders AW (2019) Timeception for complex action recognition. In: CVPR
- [59] Iqbal U, Garbade M, Gall J (2017) Pose for action-action for pose. In: International Conference on Automatic Face & Gesture Recognition, IEEE, pp 438–445
- [60] Ji Y, Yang Y, Shen HT, et al (2021) View-invariant action recognition via unsupervised attention transfer (uant). *Pattern Recognition* 113:107807
- [61] Jing L, Tian Y (2020) Self-supervised visual feature learning with deep neural networks: A survey. *PAMI*
- [62] Johnander J, Danelljan M, Brissman E, et al (2019) A generative appearance model for end-to-end video object segmentation. In: CVPR, pp 8953–8962
- [63] Kapoor R, Sharma D, Gulati T (2021) State of the art content based image retrieval techniques using deep learning: a survey. *Multimedia Tools and Applications* 80(19):29561–29583
- [64] Karbalaie A, Abtahi F, Sjöström M (2022) Event detection in surveillance videos: a review. *Multimedia Tools and Applications* 81(24):35463–35501
- [65] Karpathy A, Toderici G, Shetty S, et al (2014) Large-scale video classification with convolutional neural networks. In: CVPR
- [66] Ke L, Tai YW, Tang CK (2021) Deep occlusion-aware instance segmentation with overlapping bilayers. In: CVPR, pp 4019–4028
- [67] Ke L, Tai YW, Tang CK (2021) Occlusion-aware video object inpainting. In: ICCV, pp 14468–14478

- [68] Ke Q, Bennamoun M, An S, et al (2017) A new representation of skeleton sequences for 3d action recognition. In: CVPR, pp 3288–3297
- [69] Kim J, Li G, Yun I, et al (2021) Weakly-supervised temporal attention 3d network for human action recognition. *Pattern Recognition* 119:108068
- [70] Kim TS, Reiter A (2017) Interpretable 3d human action analysis with temporal convolutional networks. In: CVPR Workshop, IEEE, pp 1623–1631
- [71] Kniaz VV, Knyaz VA, Hladuvka J, et al (2018) Thermalgan: Multimodal color-to-thermal image translation for person re-identification in multispectral dataset. In: ECCV Workshops, pp 0–0
- [72] Kong Y, Tao Z, Fu Y (2017) Deep sequential context networks for action prediction. In: CVPR, pp 1473–1481
- [73] Kong Y, Tao Z, Fu Y (2018) Adversarial action prediction networks. *PAMI* 42(3):539–553
- [74] Korbar B, Tran D, Torresani L (2019) Scsampller: Sampling salient clips from video for efficient action recognition. In: ICCV, pp 6232–6242
- [75] Li B, Dai Y, Cheng X, et al (2017) Skeleton based action recognition using translation-scale invariant image mapping and multi-scale deep cnn. In: International Conference on Multimedia & Expo Workshops (ICMEW), IEEE, pp 601–604
- [76] Li B, Li X, Zhang Z, et al (2019) Spatio-temporal graph routing for skeleton-based action recognition. In: AAAI, pp 8561–8568
- [77] Li C, Zhong Q, Xie D, et al (2017) Skeleton-based action recognition with convolutional neural networks. In: International Conference on Multimedia & Expo Workshops, IEEE, pp 597–600
- [78] Li J, Liu X, Zhang W, et al (2020) Spatio-temporal attention networks for action recognition and detection. *IEEE Transactions on Multimedia* 22(11):2990–3001
- [79] Li L, Zheng W, Zhang Z, et al (2018) Skeleton-based relational modeling for action recognition. *arXiv preprint arXiv:180502556* 1(2):3
- [80] Li M, Chen S, Chen X, et al (2019) Actional-structural graph convolutional networks for skeleton-based action recognition. In: CVPR, pp 3595–3603
- [81] Li M, Hu L, Xiong Z, et al (2022) Recurrent dynamic embedding for video object segmentation. In: CVPR, pp 1332–1341
- [82] Li S, Bak S, Carr P, et al (2018) Diversity regularized spatiotemporal attention for video-based person re-identification. In: CVPR

- [83] Li S, Jiang T, Huang T, et al (2020) Global co-occurrence feature learning and active coordinate system conversion for skeleton-based action recognition. In: WACV, pp 586–594
- [84] Li X, Liu C, Shuai B, et al (2022) Nuta: Non-uniform temporal aggregation for action recognition. In: WACV, pp 3683–3692
- [85] Li Y, Li Y, Vasconcelos N (2018) Resound: Towards action recognition without representation bias. In: ECCV, pp 513–528
- [86] Li Y, Yang M, Zhang Z (2018) A survey of multi-view representation learning. *Transactions on knowledge and data engineering* 31(10):1863–1883
- [87] Li Y, Xia R, Liu X (2020) Learning shape and motion representations for view invariant skeleton-based action recognition. *Pattern Recognition* 103:107293
- [88] Li Y, He J, Zhang T, et al (2021) Diverse part discovery: Occluded person re-identification with part-aware transformer. In: CVPR, pp 2898–2907
- [89] Li Z, Gavriluk K, Gavves E, et al (2018) Videolstm convolves, attends and flows for action recognition. *Computer Vision and Image Understanding* 166:41–50
- [90] Liang J, Jiang L, Niebles JC, et al (2019) Peeking into the future: Predicting future person activities and locations in videos. In: CVPR, pp 5725–5734
- [91] Liang W, Zhu Y, Zhu SC (2018) Tracking occluded objects and recovering incomplete trajectories by reasoning about containment relations and human actions. In: AAAI
- [92] Liang Y, Li X, Jafari N, et al (2020) Video object segmentation with adaptive feature bank and uncertain-region refinement. *NIPS* 33:3430–3441
- [93] Lin H, Qi X, Jia J (2019) Agss-vos: Attention guided single-shot video object segmentation. In: ICCV, pp 3949–3957
- [94] Lin J, Gan C, Han S (2019) Tsm: Temporal shift module for efficient video understanding. In: ICCV, pp 7083–7093
- [95] Lin S, Xie H, Wang B, et al (2022) Knowledge distillation via the target-aware transformer. In: CVPR, pp 10915–10924
- [96] Lin Z, Yang T, Li M, et al (2022) Swem: Towards real-time video object segmentation with sequential weighted expectation-maximization. In: CVPR, pp 1362–1372
- [97] Liu D, Cui Y, Chen Y, et al (2020) Video object detection for autonomous driving: Motion-aid feature calibration. *Neurocomputing* 409:1–11

- [98] Liu D, Cui Y, Tan W, et al (2021) Sg-net: Spatial granularity network for one-stage video instance segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 9816–9825
- [99] Liu J, Shahroudy A, Xu D, et al (2016) Spatio-temporal lstm with trust gates for 3d human action recognition. In: ECCV, Springer, pp 816–833
- [100] Liu J, Akhtar N, Mian A (2017) Viewpoint invariant rgb-d human action recognition. In: International Conference on Digital Image Computing: Techniques and Applications, IEEE, pp 1–8
- [101] Liu J, Wang G, Duan LY, et al (2017) Skeleton-based human action recognition with global context-aware attention lstm networks. TIP 27(4):1586–1599
- [102] Liu J, Wang G, Hu P, et al (2017) Global context-aware attention lstm networks for 3d action recognition. In: CVPR, pp 1647–1656
- [103] Liu M, Yuan J (2018) Recognizing human actions as the evolution of pose estimation maps. In: CVPR, pp 1159–1168
- [104] Liu M, Liu H, Chen C (2017) Enhanced skeleton visualization for view invariant human action recognition. Pattern Recognition 68:346–362
- [105] Liu Y, Wang K, Li G, et al (2021) Semantics-aware adaptive knowledge distillation for sensor-to-vision action recognition. TIP 30:5573–5588
- [106] Liu Z, Zhang H, Chen Z, et al (2020) Disentangling and unifying graph convolutions for skeleton-based action recognition. In: CVPR
- [107] Liu Z, Ning J, Cao Y, et al (2022) Video swin transformer. In: CVPR, pp 3202–3211
- [108] Lu Y, Wang Q, Ma S, et al (2023) Transflow: Transformer as flow learner. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 18063–18073
- [109] Luo C, Yuille AL (2019) Grouped spatial-temporal aggregation for efficient action recognition. In: ICCV, pp 5512–5521
- [110] Luvizon DC, Picard D, Tabia H (2020) Multi-task deep learning for real-time 3d human pose estimation and action recognition. PAMI 43(8):2752–2764
- [111] Lv Z, Ota K, Lloret J, et al (2022) Complexity problems handled by advanced computer simulation technology in smart cities 2021
- [112] Ma J, Jiang X, Fan A, et al (2021) Image matching from handcrafted to deep features: A survey. IJCV 129(1):23–79

- [113] Meng Y, Lin CC, Panda R, et al (2020) Ar-net: Adaptive frame resolution for efficient action recognition. In: ECCV, Springer, pp 86–104
- [114] Minaee S, Boykov YY, Porikli F, et al (2021) Image segmentation using deep learning: A survey. PAMI
- [115] Neimark D, Bar O, Zohar M, et al (2021) Video transformer network. In: ICCV, pp 3163–3172
- [116] Oh SW, Lee JY, Xu N, et al (2019) Fast user-guided video object segmentation by interaction-and-propagation networks. In: CVPR, pp 5247–5256
- [117] Oh SW, Lee JY, Xu N, et al (2019) Video object segmentation using space-time memory networks. In: ICCV, pp 9226–9235
- [118] Ouyang W, Wang X (2012) A discriminative deep model for pedestrian detection with occlusion handling. In: CVPR, IEEE, pp 3258–3265
- [119] Ouyang W, Wang X (2013) Joint deep learning for pedestrian detection. In: ICCV, pp 2056–2063
- [120] Park K, Woo S, Oh SW, et al (2022) Per-clip video object segmentation. In: CVPR, pp 1352–1361
- [121] Patrick M, Campbell D, Asano Y, et al (2021) Keeping your eye on the ball: Trajectory attention in video transformers. NIPS 34:12493–12506
- [122] Peng W, Hong X, Chen H, et al (2020) Learning graph convolutional network for skeleton-based human action recognition by neural searching. In: AAAI, pp 2669–2676
- [123] Pexels (n.d.) Pexels. URL <https://www.pexels.com/>, accessed November 9, 2023
- [124] Piasco N, Sidibé D, Demonceaux C, et al (2018) A survey on visual-based localization: On the benefit of heterogeneous data. Pattern Recognition 74:90–109
- [125] Pont-Tuset J, Perazzi F, Caelles S, et al (2017) The 2017 davis challenge on video object segmentation. arXiv preprint arXiv:170400675
- [126] Qin X, Ge Y, Feng J, et al (2020) Dtmn: Deep transfer multi-metric network for rgb-d action recognition. Neurocomputing 406:127–134
- [127] Qin Z, Lu X, Nie X, et al (2023) Coarse-to-fine video instance segmentation with factorized conditional appearance flows. IEEE/CAA Journal of Automatica Sinica 10(5):1192–1208
- [128] Ren S, Liu W, Liu Y, et al (2021) Reciprocal transformations for unsupervised video object segmentation. In: CVPR, pp 15455–15464

- [129] Robinson A, Lawin FJ, Danelljan M, et al (2020) Learning fast and robust target models for video object segmentation. In: CVPR, pp 7406–7415
- [130] Seo S, Lee JY, Han B (2020) Urvos: Unified referring video object segmentation network with a large-scale benchmark. In: ECCV, Springer, pp 208–223
- [131] Shahroudy A, Liu J, Ng TT, et al (2016) Ntu rgb+ d: A large scale dataset for 3d human activity analysis. In: CVPR, pp 1010–1019
- [132] Sharma S, Kiros R, Salakhutdinov R (2015) Action recognition using visual attention. arXiv preprint arXiv:151104119
- [133] Shi L, Zhang Y, Cheng J, et al (2019) Skeleton-based action recognition with directed graph neural networks. In: CVPR, pp 7912–7921
- [134] Shi L, Zhang Y, Cheng J, et al (2019) Two-stream adaptive graph convolutional networks for skeleton-based action recognition. In: CVPR
- [135] Shi L, Zhang Y, Cheng J, et al (2020) Decoupled spatial-temporal attention network for skeleton-based action-gesture recognition. In: Proceedings of the Asian Conference on Computer Vision
- [136] Shi L, Zhang Y, Cheng J, et al (2020) Skeleton-based action recognition with multi-stream adaptive graph convolutional networks. TIP 29:9532–9545
- [137] Shou Z, Chan J, Zareian A, et al (2017) Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. In: CVPR
- [138] Si C, Chen W, Wang W, et al (2019) An attention enhanced graph convolutional lstm network for skeleton-based action recognition. In: CVPR, pp 1227–1236
- [139] Simonyan K, Zisserman A (2014) Two-stream convolutional networks for action recognition in videos. arXiv preprint arXiv:14062199
- [140] Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:14091556
- [141] Song L, Yu G, Yuan J, et al (2021) Human pose estimation and its application to action recognition: A survey. Journal of Visual Communication and Image Representation p 103055
- [142] Song YF, Zhang Z, Wang L (2019) Richly activated graph convolutional network for action recognition with incomplete skeletons. In: ICIP, IEEE, pp 1–5
- [143] Soomro K, Zamir AR, Shah M (2012) Ucf101: A dataset of 101 human actions classes from videos in the wild. arXiv preprint arXiv:12120402

- [144] de Souza Reis E, Seewald LA, Antunes RS, et al (2021) Monocular multi-person pose estimation: A survey. *Pattern Recognition* p 108046
- [145] Su L, Hu C, Li G, et al (2020) Msaf: Multimodal split attention fusion. *arXiv preprint arXiv:201207175*
- [146] Sudhakaran S, Escalera S, Lanz O (2020) Gate-shift networks for video action recognition. In: *CVPR*, pp 1102–1111
- [147] Sun M, Xiao J, Lim EG, et al (2020) Fast template matching and update for video object tracking and segmentation. In: *CVPR*, pp 10791–10799
- [148] Thakkar K, Narayanan P (2018) Part-based graph convolutional network for action recognition. *arXiv preprint arXiv:180904983*
- [149] Tian Y, Luo P, Wang X, et al (2015) Deep learning strong parts for pedestrian detection. In: *ICCV*, pp 1904–1912
- [150] Tran A, Cheong LF (2017) Two-stream flow-guided convolutional attention networks for action recognition. In: *ICCV Workshops*, pp 3110–3119
- [151] Tran D, Bourdev L, Fergus R, et al (2015) Learning spatiotemporal features with 3d convolutional networks. In: *ICCV*, pp 4489–4497
- [152] Tran D, Wang H, Torresani L, et al (2019) Video classification with channel-separated convolutional networks. In: *ICCV*, pp 5552–5561
- [153] Truong TD, Bui QH, Duong CN, et al (2022) Direcformer: A directed attention in transformer approach to robust action recognition. In: *CVPR*, pp 20030–20040
- [154] Ullah A, Muhammad K, Hussain T, et al (2021) Conflux lstms network: A novel approach for multi-view action recognition. *Neurocomputing* 435:321–329
- [155] Vaswani A, Shazeer N, Parmar N, et al (2017) Attention is all you need. *NIPS* 30
- [156] Veeriah V, Zhuang N, Qi GJ (2015) Differential recurrent neural networks for action recognition. In: *ICCV*, pp 4041–4049
- [157] Ventura C, Bellver M, Girbau A, et al (2019) Rvos: End-to-end recurrent network for video object segmentation. In: *CVPR*, pp 5277–5286
- [158] Voigtlaender P, Chai Y, Schroff F, et al (2019) Feelvos: Fast end-to-end embedding learning for video object segmentation. In: *CVPR*, pp 9481–9490
- [159] Wang H, Wang L (2017) Modeling temporal dynamics and spatial configurations of actions using two-stream recurrent neural networks. In: *CVPR*, pp 499–508

- [160] Wang L, Xiong Y, Wang Z, et al (2015) Towards good practices for very deep two-stream convnets. arXiv preprint arXiv:150702159
- [161] Wang L, Xiong Y, Wang Z, et al (2016) Temporal segment networks: Towards good practices for deep action recognition. In: ECCV, Springer, pp 20–36
- [162] Wang L, Tong Z, Ji B, et al (2021) Tdn: Temporal difference networks for efficient action recognition. In: CVPR, pp 1895–1904
- [163] Wang M, Ni B, Yang X (2020) Learning multi-view interactional skeleton graph for action recognition. PAMI
- [164] Wang P, Li Z, Hou Y, et al (2016) Action recognition based on joint trajectory maps using convolutional neural networks. In: Proceedings of the 24th ACM international conference on Multimedia, pp 102–106
- [165] Wang P, Li W, Gao Z, et al (2017) Scene flow to action map: A new representation for rgb-d based action recognition with convolutional neural networks. In: CVPR
- [166] Wang P, Wang S, Gao Z, et al (2017) Structured images for rgb-d action recognition. In: ICCV Workshops
- [167] Wang X, Zheng S, Yang R, et al (2022) Pedestrian attribute recognition: A survey. Pattern Recognition 121:108220. <https://doi.org/https://doi.org/10.1016/j.patcog.2021.108220>
- [168] Wang Z, Xu J, Liu L, et al (2019) Ranet: Ranking attention network for fast video object segmentation. In: ICCV, pp 3978–3987
- [169] Wen YH, Gao L, Fu H, et al (2019) Graph cnns with motif and variable temporal block for skeleton-based action recognition. In: AAAI, pp 8989–8996
- [170] Wu C, Wu XJ, Kittler J (2019) Spatial residual layer and dense connection block enhanced spatial temporal graph convolutional network for skeleton-based action recognition. In: ICCV workshops, pp 0–0
- [171] Wu D, Dong X, Shao L, et al (2022) Multi-level representation learning with semantic alignment for referring video object segmentation. In: CVPR, pp 4996–5005
- [172] Wu J, Jiang Y, Sun P, et al (2022) Language as queries for referring video object segmentation. In: CVPR, pp 4974–4984
- [173] Wu J, Yarram S, Liang H, et al (2022) Efficient video instance segmentation via tracklet query and proposal. In: CVPR

- [174] Wu W, He D, Tan X, et al (2019) Multi-agent reinforcement learning based frame sampling for effective untrimmed video recognition. In: ICCV, pp 6222–6231
- [175] Xie H, Yao H, Zhou S, et al (2021) Efficient regional memory network for video object segmentation. In: CVPR, pp 1286–1295
- [176] Xie S, Sun C, Huang J, et al (2018) Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In: ECCV, pp 305–321
- [177] Xu C, Govindarajan LN, Zhang Y, et al (2017) Lie-x: Depth image based articulated object pose estimation, tracking, and action recognition on lie groups. IJCV 123(3):454–478
- [178] Xu J, Zhao R, Zhu F, et al (2018) Attention-aware compositional network for person re-identification. In: CVPR, pp 2119–2128
- [179] Xu K, Yao A (2022) Accelerating video object segmentation with compressed video. In: CVPR, pp 1342–1351
- [180] Xu K, Wen L, Li G, et al (2019) Spatiotemporal cnn for video object segmentation. In: CVPR, pp 1379–1388
- [181] Xu M, Gao M, Chen YT, et al (2019) Temporal recurrent networks for online action detection. In: ICCV, pp 5532–5541
- [182] Xu N, Yang L, Fan Y, et al (2018) Youtube-vos: A large-scale video object segmentation benchmark. arXiv preprint arXiv:180903327
- [183] Xu S, Cheng Y, Gu K, et al (2017) Jointly attentive spatial-temporal pooling networks for video-based person re-identification. In: ICCV, pp 4733–4742
- [184] Yan A, Wang Y, Li Z, et al (2019) Pa3d: Pose-action 3d machine for video recognition. In: CVPR
- [185] Yan A, Wang Y, Li Z, et al (2019) Pa3d: Pose-action 3d machine for video recognition. In: CVPR, pp 7922–7931
- [186] Yan L, Wang Q, Cui Y, et al (2022) Gl-rg: Global-local representation granularity for video captioning. arXiv preprint arXiv:220510706
- [187] Yan S, Xiong Y, Lin D (2018) Spatial temporal graph convolutional networks for skeleton-based action recognition. In: AAAI
- [188] Yang H, Yuan C, Li B, et al (2019) Asymmetric 3d convolutional neural networks for action recognition. Pattern Recognition 85:1–12
- [189] Yang H, Yan D, Zhang L, et al (2021) Feedback graph convolutional network for skeleton-based action recognition. TIP 31:164–175

- [190] Yang J, Dong X, Liu L, et al (2022) Recurring the transformer for video action recognition. In: CVPR, pp 14063–14073
- [191] Yang L, Fan Y, Xu N (2019) Video instance segmentation. In: CVPR, pp 5188–5197
- [192] Yu F, Koltun V (2015) Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:151107122
- [193] Zhang D, Dai X, Wang YF (2018) Dynamic temporal pyramid network: A closer look at multi-scale modeling for activity detection. In: Asian Conference on Computer Vision, Springer, pp 712–728
- [194] Zhang K, Zhao Z, Liu D, et al (2021) Deep transport network for unsupervised video object segmentation. In: ICCV, pp 8781–8790
- [195] Zhang L, Lin Z, Zhang J, et al (2019) Fast video object segmentation via dynamic targeting network. In: ICCV, pp 5582–5591
- [196] Zhang P, Lan C, Xing J, et al (2017) View adaptive recurrent neural networks for high performance human action recognition from skeleton data. In: ICCV, pp 2117–2126
- [197] Zhang R, Li J, Sun H, et al (2019) Scan: Self-and-collaborative attention network for video person re-identification. TIP 28(10):4870–4882
- [198] Zhang S, Yang J, Schiele B (2018) Occluded pedestrian detection through guided attention in cnns. In: CVPR, pp 6995–7003
- [199] Zhang Y, Borse S, Cai H, et al (2022) Perceptual consistency in video segmentation. In: WACV, pp 2564–2573
- [200] Zhao H, Wildes RP (2019) Spatiotemporal feature residual propagation for action prediction. In: ICCV, pp 7003–7012
- [201] Zhao L, Wang Y, Zhao J, et al (2021) Learning view-disentangled human pose representation by contrastive cross-view mutual information maximization. In: CVPR, pp 12793–12802
- [202] Zheng Z, An G, Wu D, et al (2020) Global and local knowledge-aware attention network for action recognition. IEEE Transactions on Neural Networks and Learning Systems 32(1):334–347
- [203] Zhou C, Yuan J (2017) Multi-label learning of part detectors for heavily occluded pedestrian detection. In: ICCV, pp 3486–3495
- [204] Zhou Q, Sheng K, Zheng X, et al (2022) Training-free transformer architecture search. In: CVPR, pp 10894–10903

- [205] Zhou Y, Zhang H, Lee H, et al (2022) Slot-vps: Object-centric representation learning for video panoptic segmentation. In: CVPR, pp 3093–3103
- [206] Zhu D, Zhang Z, Cui P, et al (2019) Robust graph convolutional networks against adversarial attacks. In: Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining, pp 1399–1407
- [207] Zhu J, Zou W, Xu L, et al (2018) Action machine: Rethinking action recognition in trimmed videos. arXiv preprint arXiv:181205770
- [208] Zolfaghari M, Singh K, Brox T (2018) Eco: Efficient convolutional network for online video understanding. In: ECCV, pp 695–712
- [209] Zolfaghari M, Zhu Y, Gehler P, et al (2021) Crossclr: Cross-modal contrastive learning for multi-modal video representations. In: ICCV, pp 1450–1459
- [210] Zong M, Wang R, Chen X, et al (2021) Motion saliency based multi-stream multiplier resnets for action recognition. *Image and Vision Computing* 107:104108