# A Review of Video Object Detection: Datasets, Metrics and Methods

**Haidi Zhu [1,2], Haoran Wei [3] 🔾, Baoqing Li [1,*] 🔾, Xiaobing Yuan [1] and Nasser Kehtarnavaz [3]**

1   Science and Technology on Micro-System Laboratory,
    Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences,
    Shanghai 201800, China; hdzhu@mail.sim.ac.cn (H.Z.); sinowsn@mail.sim.ac.cn (X.Y.)
2   University of Chinese Academy of Sciences, Beijing 100049, China
3   Department of Electrical and Computer Engineering, University of Texas at Dallas, Richardson, TX 75080,
    USA; Haoran.Wei@utdallas.edu (H.W.); kehtar@utdallas.edu (N.K.)
*   Correspondence: sinoiot@mail.sim.ac.cn

check for
updates

**Abstract:** Although there are well established object detection methods based on static images, their application to video data on a frame by frame basis faces two shortcomings: (i) lack of computational efficiency due to redundancy across image frames or by not using a temporal and spatial correlation of features across image frames, and (ii) lack of robustness to real-world conditions such as motion blur and occlusion. Since the introduction of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2015, a growing number of methods have appeared in the literature on video object detection, many of which have utilized deep learning models. The aim of this paper is to provide a review of these papers on video object detection. An overview of the existing datasets for video object detection together with commonly used evaluation metrics is first presented. Video object detection methods are then categorized and a description of each of them is stated. Two comparison tables are provided to see their differences in terms of both accuracy and computational efficiency. Finally, some future trends in video object detection to address the challenges involved are noted.

## 1. Introduction

Video object detection involves detecting objects using video data as compared to conventional object detection using static images. Two applications that have played a major role in the growth of video object detection are autonomous driving [1,2] and video surveillance [3,4]. In 2015, video object detection became a new task of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC2015) [5]. With the help of ILSVRC2015, studies in video object detection have further increased.

Earlier attempts in video object detection involved performing object detection on each image frame. In general, object detection approaches can be grouped into two major categories: (1) one-stage detectors and (2) two-stage detectors. One-stage detectors (e.g., [6–12]) are often more computationally efficient than two-stage detectors (e.g., [13–21]). However, two-stage detectors are shown to produce higher accuracies compared to one-stage detectors.

However, using object detection on each image frame does not take into consideration the following attributes in video data: (1) Since there exist both spatial and temporal correlations between image frames, there are feature extraction redundancies between adjacent frames. Detecting features in each frame leads to computational inefficiency. (2) In a long video stream, some frames may have poor

quality due to motion blur, video defocus, occlusion, and pose changes [22]. Detecting objects from poor quality frames leads to low accuracies. Video object detection approaches attempt to address the above challenges. Some approaches make use of the spatial-temporal information to improve accuracy, such as fusing features on different levels, e.g., [22–25]. Some other approaches focus on reducing information redundancy and improving detection efficiency, e.g., [26–28].

Initially, video object detection approaches have relied on handcrafted features, e.g., [29–42]. With the rapid development of deep learning and convolutional neural networks, deep learning models have been shown to be more effective than conventional approaches for various tasks in computer vision [43–50], speech processing [51–55], and multi-modality signal processing [56–61]. A number of deep learning-based video object detection approaches were developed after the ILSVRC2015 challenge. The training is normally done offline. The testing phase on modern GPUs even of complex networks has been shown to meet the 30 frames per sec rate of video, e.g., [26], allowing the real-time deployment of networks.

The great value of video object detection approaches is further presented in some specific applications. For example, hand segmentation [62,63] is well realized with the help of the optical flow to enhance the feature maps as per the video object detection method [28]. Human pose estimation in videos [64] is another successful application, which draws lessons from [22,28] to solve the motion blur, occlusion and other specific challenges occurring in videos. Furthermore, instance-level human parsing [65] starts from the similar approaches. Mutual assistance of tracking and detection [26] is well employed in multiple people tracking [66].

Deep learning-based video object detection approaches can be divided into flow based [22,27,28,67–69], LSTM (Long Short Term Memory)-based [70–73], attention-based [25,74–77], tracking-based [26,78–82] and other methods [36,83–90]. A review of these approaches is provided in this paper.

Section 2 covers the existing datasets and evaluation metrics for video object detection. Then, in Section 3, the existing video object detection approaches are described. The accuracy and processing time of these approaches are compared in Section 4. Section 5 mentions the future trends or needs relating to video object detection. Finally, the conclusion is stated in Section 6.

## 2. Datasets and Evaluation Metrics

### 2.1. Datasets

Many datasets have been provided for specific applications [91–93]. For video object detection, the most commonly used dataset is the ImageNet VID dataset [5], which is a prevalent benchmark for video object detection. The dataset is split into a training set and a validation set, containing 3862 video snippets and 555 video snippets, respectively. The video streams are annotated on each frame at the frame rate of 25 or 30 fps. In addition, this dataset contains 30 object categories, which are a subset of the categories in the ImageNet DET dataset [93].

In the ImageNet VID dataset, the number of objects in each frame is small compared with the datasets used for static image object detection such as COCO [92]. Though the ImageNet VID dataset is widely used, it has limitations in fully reflecting the effect of various video object detection methods. In [94], a large-scale dataset named YouTube-BoundingBoxes (YT-BB) was provided, which is human-annotated at one frame per s on video snippets from YouTube with high accuracy classification labels and tight bounding boxes. YT-BB contains approximately 380,000 video segments with 5.6 million bounding boxes of 23 object categories, which is a subset of the COCO label set. However, the dataset contains only 23 object categories and the image quality is relatively low due to its collection by hand-held mobile phones.

In 2018, a dataset named EPIC KITCHENS was provided in [95], which consists of 32 different kitchens in 4 cities with 11,500,000 frames containing 454,158 bounding boxes spanning 290 classes. However, its kitchen scenario poses limitations on performing generic video object detection. Moreover,

there exist the following other datasets that reflect specific applications: the DAVIS dataset [96] for object segmentation, CDnet2014 [97] for moving object detection, VOT [98] and MOT [99] for object tracking, Sports-1M data set [100] with segment-level annotations, HMDB-51 data set [101] with segment-level annotations for various human action categories, TRECVID [102] for video retrieval and indexing, the Caltech Pedestrian Detection data set [103] for pedestrian detection, and the PASCAL VOC dataset [104,105] for object detection. In addition, some works based on semi-supervised or unsupervised methods have been considered in [106–109].

For video object detection with classification labels and tight bounding boxes annotation, currently there exists no public domain dataset offering dense annotations for various complex scenes. To enable the advancement of video object detection, more effort is thus needed to establish comprehensive datasets.

## 2.2. Evaluation Metrics

The metric mean Average Precision (mAP) is extensively used in conventional object detection, which provides a performance evaluation in terms of regression and classification accuracies [9–15,17]. The evaluation metric mAP represents the mean Average Precision. The definition is the mean of the Average Precision of each category. As per the PASCAL Visual Object Classes Challenge 2012 (VOC2012) Development Kit, it is computed as follows:

(1)  The Precision/Recall curve is obtained first. For the Recall (r), the Precision is set to the maximum Precision achieved for any Recall $r' \geq r$.
(2)  The area under the Precision/Recall curve is considered to be the Average Precision (AP). The mean of AP in each category is mAP.

Prior to 2010, AP used to be computed by sampling the curve at a set of uniformly spaced Recall $(0, 0.1, 0.2, \ldots, 1)$ and then computing the average of the corresponding Precision value. More specifically, Recall $= \frac{TP}{TP+FN}$ and Precison $= \frac{TP}{TP+FP}$, where the definitions of TP, FP and FN appear in Table 1.

**Table 1.** Definitions of TP, FP, TN and FN.

|                      | Label is True.      | Label is False.     |
| -------------------- | ------------------- | ------------------- |
| Prediction is true.  | True positive (TP)  | False positive (FP) |
| Prediction is false. | False negative (FN) | True negative (TN)  |

When the IoU is larger than a set threshold, the prediction is true. That is, $\text{IoU} = \frac{\text{area}(B_{gt} \cap B_p)}{\text{area}(B_{gt} \cup B_p)}$, where $B_{gt}$ and $B_p$ indicate the ground truth and prediction box, respectively. More details are stated in the following example.

The detection results of one category are presented in Table 2, and the number of the objects is 3, which means that TP + FN = 3. Confidence represents the confidence level of the prediction boxes. The definition of confidence score is stated in Equation (1) with $P_r$ denoting precision and $\text{IoU}_{\text{pred}}^{\text{truth}}$ the confidence level of the box surrounding the entire object,

$$\text{confidence score} = P_r(\text{object}) \times \text{IoU}_{\text{pred}}^{\text{truth}},$$
$$P_r(\text{object}) \in [0, 1]. \tag{1}$$

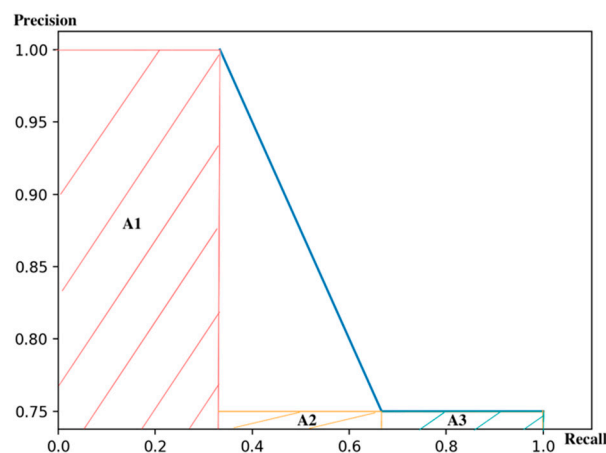**Table 2.** Example detection results of one category.

| Image Detection | Confidence Score | TP or FP |
|:---:|:---:|:---:|
| 1 | 0.92 | TP |
| 2 | 0.83 | TP |
| 3 | 0.85 | FP |
| 4 | 0.75 | TP |
| 5 | 0.72 | FP |

The detection results are ranked according to the confidence score, which are shown in Table 3. The Precision and Recall are computed by the Equation noted above. The Precision/Recall curve of this category is shown in Figure 1. As a result, $AP = A1 + A2 + A3 = \left(\frac{1}{3} \times 1\right) + \left(\frac{2}{3} - \frac{1}{3}\right) \times \frac{3}{4} + \left(1 - \frac{2}{3}\right) \times \frac{3}{4} = 83.33\%$ and mAP is the mean of the AP in each category.

**Table 3.** Ranked detection results of one category according to the confidence score.

| Image Detection | Confidence Score | TP or FP | Precision | Recall | Maximum Precision for any Recall $r' \geq r$ |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | 0.92 | TP | 1/1 | 1/3 | |
| 3 | 0.85 | FP | 1/2 | 1/3 | 1 |
| 2 | 0.83 | TP | 2/3 | 2/3 | 3/4 |
| 4 | 0.75 | TP | 3/4 | 3/3 | |
| 5 | 0.72 | FP | 3/5 | 3/3 | 3/4 |



**Figure 1.** The precision/recall curve.

For video object detection, mAP is also directly used as an evaluation metric in [22,25,28,72,74]. Based on the object speed, it is labeled as mAP (slow), mAP (medium) and mAP (fast) [22]. This is done using the average score of IoU (Intersection over Union) of a current frame and 10 frames ahead and past as follows: slow (score > 0.9), medium (score $\in$ [0.7, 0.9]) and fast (score < 0.7).

In [110], it was pointed out that performance cannot be sufficiently evaluated using only Average Precision (AP), since the temporal nature of video snippets do not get captured by it. In the same paper, a new metric named Average Delay (AD) was introduced based on the number of frames taken to detect an object starting from the frame it first appears in. A subset of the ImageNet VID dataset, named ImageNet VIDT, was considered to verify the effectiveness of AD. It has been reported that most methods with higher ADs still had good APs or good average detection accuracies. However, higher ADs also mean that the detection delay is large. In other words, the number of frames from the frame that the object first appears in is large. If only using AP as the metric to evaluate the performance of different methods, it becomes challenging to reflect the AD (the number of frames taken to detect an object starting from the frame it first appears). As a result, AP is not sufficient to reflect the temporal characteristics of video object detectors and the metric AD provides a complementary performance indicator.

## 3. Video Object Detection Methods

For video object detection, in order to make full use of the video characteristics, different methods are considered to capture the temporal–spatial relationship. Some papers have considered the traditional methods [29–42]. These papers heavily rely on the manual design leading to the shortcomings of low accuracy and the lack of robustness to noise sources. More recently, deep learning solutions have attempted to overcome these shortcomings. As shown in Figure 2, based on the utilization of the temporal information and the aggregation of features extracted from video snippets, video object detectors can be divided into flow-based [22,27,28,67–69], LSTM-based [70–73], attention-based [25,74–77], tracking-based [26,78–82] and other methods [36,83–90]. These methods are described in more detail below.
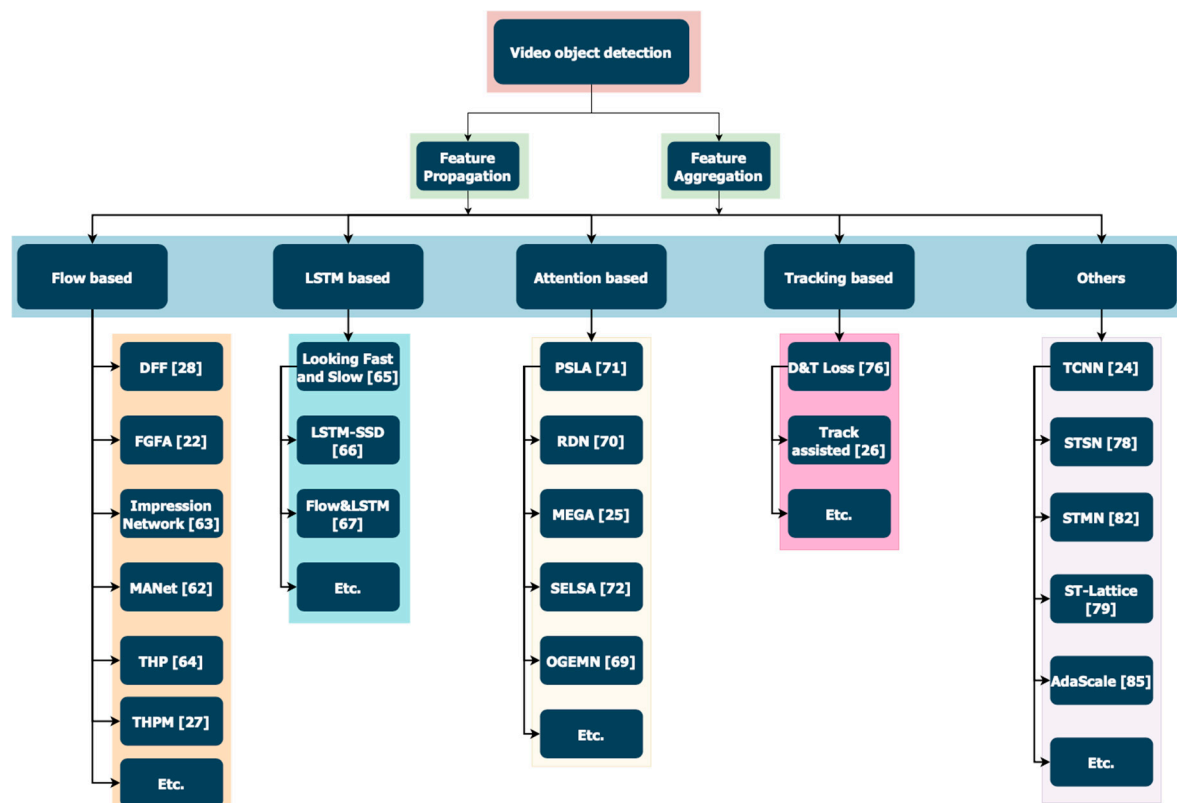


**Figure 2.** Categories of video object detection methods.

### 3.1. Flow-Based

Flow-based methods use optical flow in two ways. In order to save computation, in the first way as discussed in [28] (DFF (DFF is the acronym standing for Deep Feature Flow for Video Recognition in reference [28]. Similarly, other acronyms appear in the references indicated)), optical flow is used to propagate features from key frames to non-key frames. In the second way, as discussed in [22] (FGFA), optical flow is used to make use of the temporal–spatial information between adjacent frames to enhance the features of each frame. In the second way, higher detection accuracies but lower speeds are reported. As a result, attempts were made to combine both of these ways in [68] (Impression Network), [69] (THP) and [27] (THPM). To obtain the difference between adjacent frames and utilize the temporal–spatial information at the pixel level, an optical flow algorithm was proposed in [29]. In [111], the optical flow estimation was achieved by using the deep learning model of FlowNet.

For video object detection, it is challenging to apply the state-of-the-art object detection approaches for still images directly to each image frame in video data for the reasons stated earlier. Therefore, based on FlowNet, the DFF method was proposed in [28] to address these shortcomings: (i) computation time

of feature map extraction for each frame in video, (ii) similarity of features obtained on two adjacent frames, (iii) propagation of feature maps from one frame to another. In [28], a convolutional neural sub-network, ResNet-101, was employed to extract the feature map on sparse key frames. Features on non-key frames were obtained by warping the feature map on key frames with the flow field generated by FlowNet [111] instead of getting extracted by ResNet-101. The framework is shown in Figure 3. This method accelerates the object detection on non-key frames. On the ImageNet VID dataset [5], DFF achieved an accuracy of 73.1% mAP with 20 fps, while the baseline accuracy on a single frame was 73.9% with 4 fps. This method significantly advanced the practical aspect of video object detection.
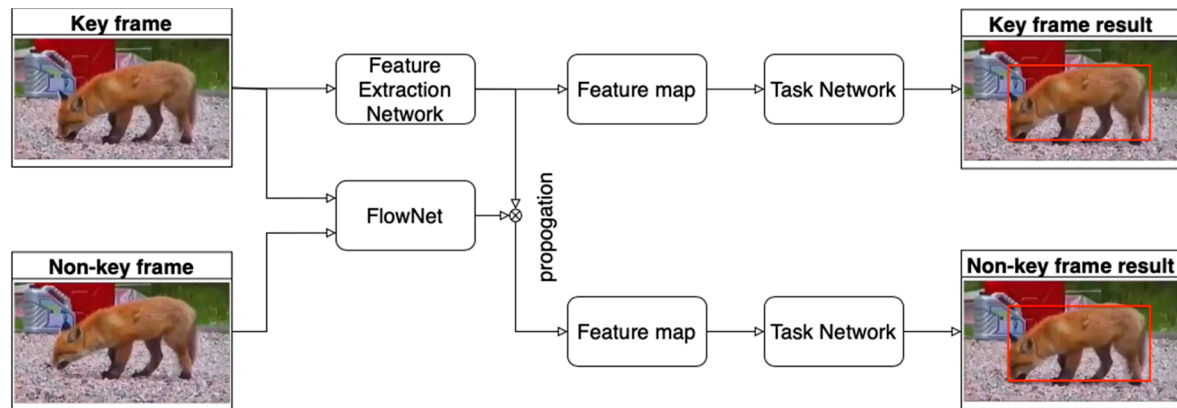


**Figure 3.** DFF (deep feature flow) framework [28].

In [22], a flow-guided feature aggregation (FGFA) method was proposed to improve the detection accuracy due to motion blur, rare poses, video defocus, etc. Feature maps were extracted on each frame in video using ResNet-101 [112]. In order to enhance the feature maps of a current frame, the feature maps of its nearby frames were warped to the current frame according to the motion information obtained by the optical flow network. The warped feature maps and extracted feature maps on the current frame were then inputted into a small sub-network to obtain a new embedding feature, which was used for a similarity measure based on the cosine similarity metric [113] to compute the weights. Next, the features were aggregated according to the weights. Finally, the aggregated feature maps were inputted into a shallow detection-specific sub-network to obtain the final detection outcome on the current frame. The framework of FGFA is shown in Figure 4. Based on the ImageNet VID dataset, FGFA achieved an accuracy of 76.3% mAP with 1.36 fps, which was higher than DFF.

Although the feature fusion method of FGFA improved the detection accuracy, it considerably increased the computation time. On the other hand, feature propagation methods showed improved computational efficiency but at the expense of reduced detection accuracy. In 2017, a so-called Impression Network [68] was developed to improve the performance in terms of both accuracy and computational speed simultaneously. Inspired by the idea that humans do not forget the previous frames when a new frame is observed, sparse key-frame features were aggregated with other key frames to improve the detection accuracy. Feature maps of non-key frames were also obtained by a feature propagation method similar to that in [28] with the assistant of a flow field. As a result, feature propagation to obtain the features of the non-key frames improved the inference computation speed. The feature aggregation method on the key frames used a small fully convolutional network to obtain the weight maps on each localization, which was different from the method in [22]. The Impression Network achieved 75.5% mAP accuracy at 20 fps on the ImageNet VID dataset.
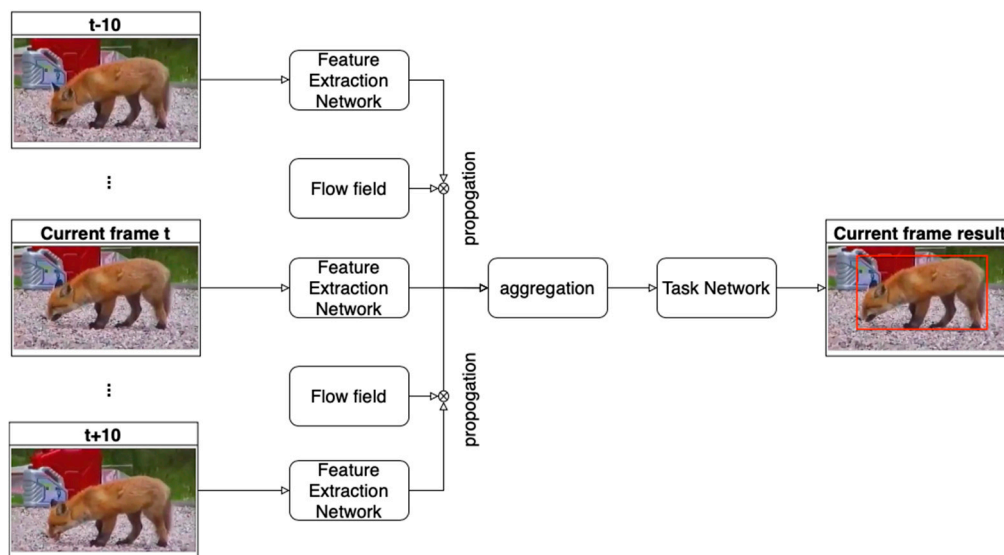
**Figure 4.** FGFA (flow-guided feature aggregation) framework [22].

Besides Impression Network, in [69], another combination method (THP) was introduced. Noting that all of the above methods utilized fixed interval key frames, this method introduced a temporally adaptive key frame scheduling to further improve the trade-off between speed and accuracy. Fixed interval key frames pose a difficulty to control the quality of key frames. With temporally adaptive key frame scheduling, the fixed interval key frames were adjusted in a dynamic manner according to the proportion of points with poor optical flow quality. If it was greater than a prescribed threshold T, it would indicate that a current frame had changed too much compared with the previous key frame. The current frame was then chosen as the new key frame and the feature maps were obtained from it.

According to the results reported in [69], the mAP accuracy was 78.6% with a runtime of 13.0 and 8.6 fps on the GPUs NVIDIA Titan X and K40 (NVIDIA, California, USA) respectively. With a different T, the mAP slightly decreased to 77.8% at faster speeds (22.9 and 15.2 fps on Titan X and K40, respectively). Compared with the winning entry [114] of the ImageNet VID challenge 2017, which was based on feature propagation [28] and aggregation [22], an mAP of 76.8% at 15.4 fps was achieved on Titan X, and a better performance in terms of both the detection accuracy and speed was obtained in [69].

Similarly, THPM [27] provided a light weight network architecture for video object detection. A light image object detector is utilized on key frames. The state-of-the-art lightweight Mobilenet [115] is utilized as the backbone network. Feature maps from key frames are propagated to non-key frame for detection by a light flow network. A flow-guided gated recurrent unit (GRU) module is provided to aggregate features effectively between key frames. On the ImageNet VID dataset, THPM achieves 60.2% mAP at speed of 25.6 fps on mobiles (e.g., HuaWei Mate 8 produced by HUAWEI TECHNOLOGIES CO., LTD, China).

### 3.2. LSTM-Based

In order to make full use of the temporal–spatial information, convolutional long short term memory (LSTM [116]) was employed to process sequential data in [117] and select important information for a long duration. The methods reported in [70] and [71] are offline LSTM-based solutions, which utilize all the frames in the video. While the method in [72] is an online solution, it only uses the current and previous frames.

In [71], a light model was proposed, which was designed to work on mobile phones and embedded devices. This method integrated SSD [9] (an efficient object detector network) with the convolutional LSTM by applying an image-based object detector to video object detection via a convolutional LSTM.

The convolutional LSTM was a modified version of the traditional LSTM encoding the temporal and spatial information.

Considering a video snippet as video frames V = {$I_0$, $I_1$, $I_2$, … $I_t$}, the model is viewed as a function $F(I_t, S_{t-1}) = (D_t, S_t)$, where $D_t$ denotes the detection outcome of the video object detector and $S_t$ represents a vector of feature maps up to the video frame $t$. Each feature map of $S_{t-1}$ is the state input to the LSTM and $S_t$ is the state output. The state unit $S_t$ of LSTM contains the temporal information. LSTM can combine the state unit with input features, adaptively adding the temporal information to the input features, and updating the state unit at the same time. In [71], it was stated that such a convolutional LSTM layers could be added to any layer of the original object detector to refine the input features of the next layer. An LSTM layer could be placed immediately after any feature map. Placing the LSTM earlier would lead to larger input volumes and much higher computational costs. In [71], the convolutional LSTM was placed only after the Conv13 layer, which was proved to be most effective through experimental analysis. This method was evaluated on the ImageNet VID 2015 dataset [5] and achieved a good performance in terms of the model size and computational efficiency (15 fps on a mobile CPU), with an accuracy comparable to those more computationally demanding single frame models.

In 2019, the method in [71] was improved in [70] in terms of inference speed. Specifically, as shown in Figure 5, due to the high temporal redundancy in the video, the model proposed in [70] contained two feature extractors: a small feature extractor and a large feature extractor. The large feature extractor with low speed was responsible for extracting the features with high accuracy, while the small feature extractor with a fast speed was responsible for extracting the features with poor accuracy. The two feature extractors were used alternately. The feature maps were aggregated using a memory mechanism with the modified convolutional LSTM layer. Then, a SSD-style [9] detector was applied to the refined features to obtain the final regression and classification outcome.
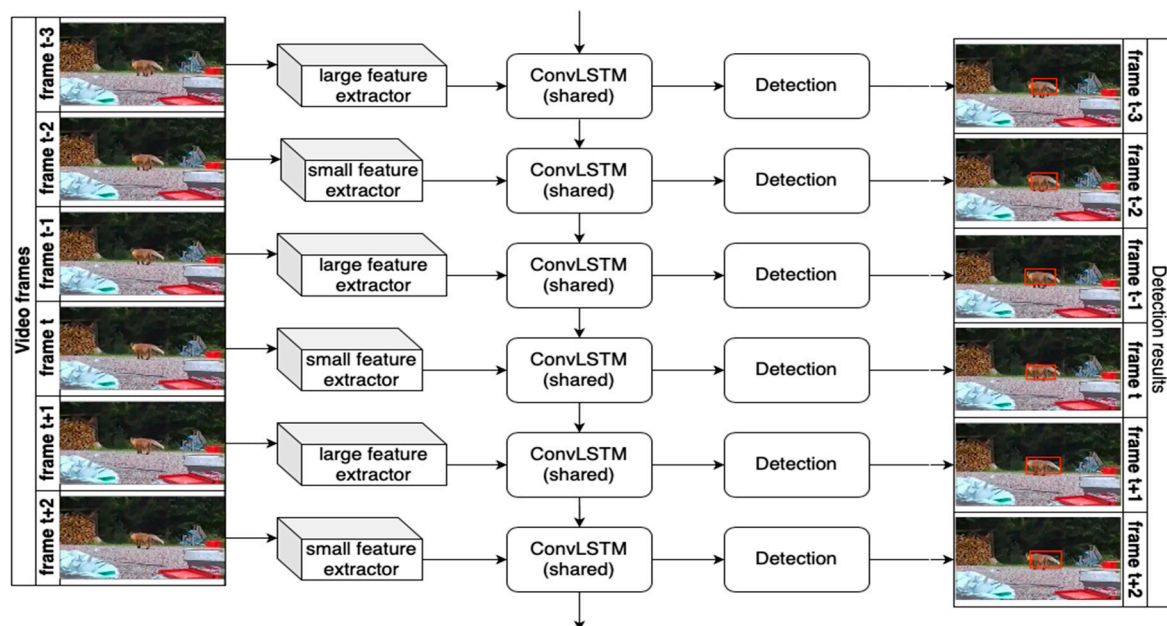


**Figure 5.** Small and large feature extractors in [70].

For the methods mentioned above, image object detectors together with a temporal context information enhancement were employed to detect objects in video. However, for online video object detection, succeeding frames cannot be utilized. In other words, non-causal video object detectors are not feasible for online applications. Noting that most video object detectors are non-causal, a causal recurrent method was proposed in [72] for online detection without using succeeding frames. In this case, the challenges in terms of occlusion and motion blur remain, which requires the use of temporal

information. For online video object detection, only the current frame and the previous frame are used. Based on the optical flow method [111], the short-term temporal information was utilized by warping the feature maps from the previous frame. However, sometimes image distortion or occlusion would last for several video frames. By using only the short-term temporal information, it was difficult to deal with these situations. The long-term temporal context information was also exploited via the convolutional LSTM, in which the feature maps of the distant preceding frame obtained from the memory function was propagated to acquire more information. The important sub-network (temporal Conv LSTM) is shown in Figure 6. Given the feature map at the time step t, the state and output from the time step t−1, the output $H_t$ and the updated state $S_t$ at the current time step t are computed as Equation (2). The long-term temporal information is stored, propagated and employed. Then, the feature maps extracted on the current frame as well as the warped feature maps and the output of the LSTM were concatenated to obtain the aggregated feature maps. Finally, the aggregated feature maps were inputted into a detection sub-network to obtain the detection outcome on the current frame. By utilizing both the short- and long-term information, this method achieved an accuracy of 75.5% mAP at a high speed on the ImageNet VID dataset, indicating a competitive performance for online detection.

$$
\begin{aligned}
FG_t &= \sigma(W_{FG} * concat(F_t, H_{t-1})), \\
I_t &= \sigma(W_I * concat(F_t, H_{t-1})), \\
O_t &= \sigma(w_O * concat(F_t, H_{t-1})), \\
C_t &= tanh(w_C * concat(F_t, H_{t-1})), \\
S_t &= (f_t \times S_{t-1}) + (I_t \times C_t), \\
H_t &= tanh(S_t) \times O_t.
\end{aligned}
\tag{2}
$$

where the $FG_t$, $I_t$, $O_t$ and $C_t$ denote the output of Forget Gate, Input Gate, Output Gate and the information branch at the time step t, respectively. Their weights are represented by $W_{FG}$, $W_I$, $W_O$ and $W_C$. $\sigma()$, $\times$, $+$, $*$ represent the activation function, element-wise multiplication, element-wise addition and 3*3 convolutions operations, respectively.
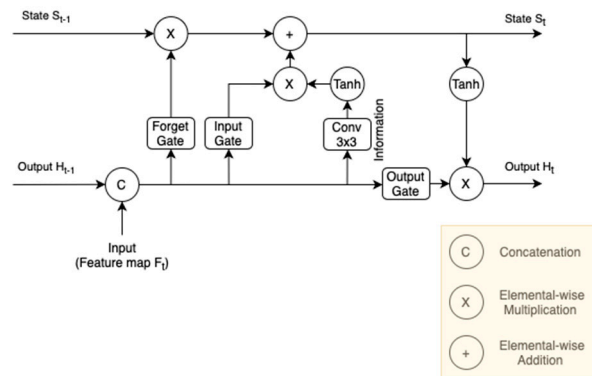


**Figure 6.** Temporal Conv LSTM architecture. State $S_{t-1}$ and Output $H_{t-1}$ are retrieved from the memory. Forget Gate, Input Gate and Output Gate operate the $3 \times 3$ convolutions followed by the activation function.

## 3.3. Attention-Related

For video object detection, it is known that exploiting the temporal context relationship is quite important. This relationship needs to be established based on a long-duration video, which requires a large amount of memory and computational resources. In order to decrease the computational resources, an attention mechanism was introduced for feature map alignment. This mechanism was first proposed for machine translation in [118,119] and was then applied to video object detection in [25,74–77].

Some methods only take the global or local temporal information into consideration. Specifically, the method RDN in [75] only makes use of the local temporal information. The methods SELSA in [77], and OGEMN in [74] only utilize the global temporal information. While the other methods of PSLA in [76], and MEGA in [25] use both the global and local temporal information.

Relation Distillation Networks (RDN) presented in [75] propagate and aggregate the feature maps based on object relationships in video. In RDN, ResNet-101 [112] and ResNeXt-101-64 × 4d [120] are utilized as the backbone to extract feature maps and object proposals are generated with the help of a Region Proposal Network (RPN) [15]. The feature maps of each proposal on the reference frame are augmented on the basis of supportive proposals. A prominent innovation in this work is to distill the relation with multi-stage reasoning consisting of a basic and an advanced stage. In the basic stage, the supportive proposals consisting of Top K proposals of a current frame and its adjacent frames are used to measure the relation feature of each reference proposal obtained on the current frame to generate refined reference proposals. In the advanced stage, supportive proposals with high objective scores are selected to generate advanced supportive proposals. Features of selected supportive proposals are aggregated with the relation against all supportive proposals. Then, such aggregated features are employed to strengthen the reference proposals obtained from the basic stage. Finally, the aggregated features of reference proposals obtained from the advanced stage are used to generate the final classification and bounding box regression. In addition, the detection box linking is used in a post-processing stage to refine the detection outcome. Evaluated on the ImageNet VID dataset, RDN achieved a detection accuracy of 81.8% and 83.2% mAP, respectively, with ResNet-101 and ResNeXt-101 for feature extraction. With linking and rescoring operations, it achieved an accuracy of 83.8% and 84.7% mAP, respectively.

A module (SELSA) was introduced in [77] to exploit the relationship between the proposals in the entire sequence level, and then related feature maps were fused for classification and regression. More specifically, the features of the proposals were extracted on different frames and then a clustering module and a transformation module were applied. The similarities of the proposals were computed across frames and the features were aggregated according to the similarities. Consequently, more robust features were generated for the final detection.

In [74], OGEMN was presented and used object-guided external memory to store the pixel and instance level features for further global aggregation. In order to improve the storage-efficiency aspect, only the features within the bounding boxes were stored for further feature aggregation.

In [25], MEGA was introduced, utilizing the global and local information inspired by how humans go about object detection in video using both global semantic information and local localization information. For situations when it was difficult to determine what the object was in the current frame, the global information was utilized to recognize a fuzzy object according to a clear object with a high similarity in another frame. When it was difficult to find out where the object was in a frame, the local localization information was used by taking the difference between adjacent frames if it was moving. More specifically, RPN was used to generate candidate proposals from those local frames (adjacent frames of current frames) and global frames. Then, a relation module was set up to aggregate the features of candidate proposals on global frames into that of local frames. This was named the global aggregation stage. With this method, the global information was integrated into the local frames. Then, features of the current frame were further augmented by the relation modules in the local aggregation stage. In order to expand the aggregation scale, an efficient module (Long Range Memory (LRM)) was designed where all the features computed in the middle were saved and utilized in a following detection. Evaluated on the ImageNet VID dataset, MEGA with ResNet-101 as backbone achieved an accuracy of 82.9% mAP. Compared with the competitor RDN, MEGA produced 1.1% improvement. Replacing ResNet-101 with ResNeXt-101 or with a stronger backbone to extract features, MEGA obtained an accuracy of 84.1% mAP. With the help of post-processing, it achieved 1.6% and 1.3% improvement with ResNet-101 and ResNeXt-101, respectively.

The method Progressive Sparse Local Attention (PSLA) was proposed in [76] to make use of the long term temporal information for enhancement on each feature cell in an attention manner. PSLA establishes correspondence by propagating features in a local region with a gradually sparser stride according to the spatial information across frames. Recursive Feature Updating (RFU) and Dense Feature Transforming (DenseFT) were also proposed based on PSLA to model the temporal relationship and enhance the features in a framework shown in Figure 7. More specifically, features were propagated in an attention manner. First, the correspondence between each feature cell in an embedding feature map of a current frame and its surrounding cells was established with a progressive sparser stride from the center to the outside of another embedding feature map of a support frame. Second, correspondence weights were used to compute the aligned feature maps. The feature maps were aggregated with the aligned features. In addition, similar to other video object detectors, the features of key frames were propagated to non-key frames. A lightweight network was then applied to extract low-level features on non-key frames and fuse them with the features propagated from key frames (DenseFT). Feature propagation was also employed between key frames, and key frame features were updated recursively by an update network (RFU). Hence, features were enriched by the temporal information with DenseFT and RFU, which were further used for detection. Based on the experimentations done in [76], an accuracy of 81.4% mAP was achieved on the ImageNet VID dataset.
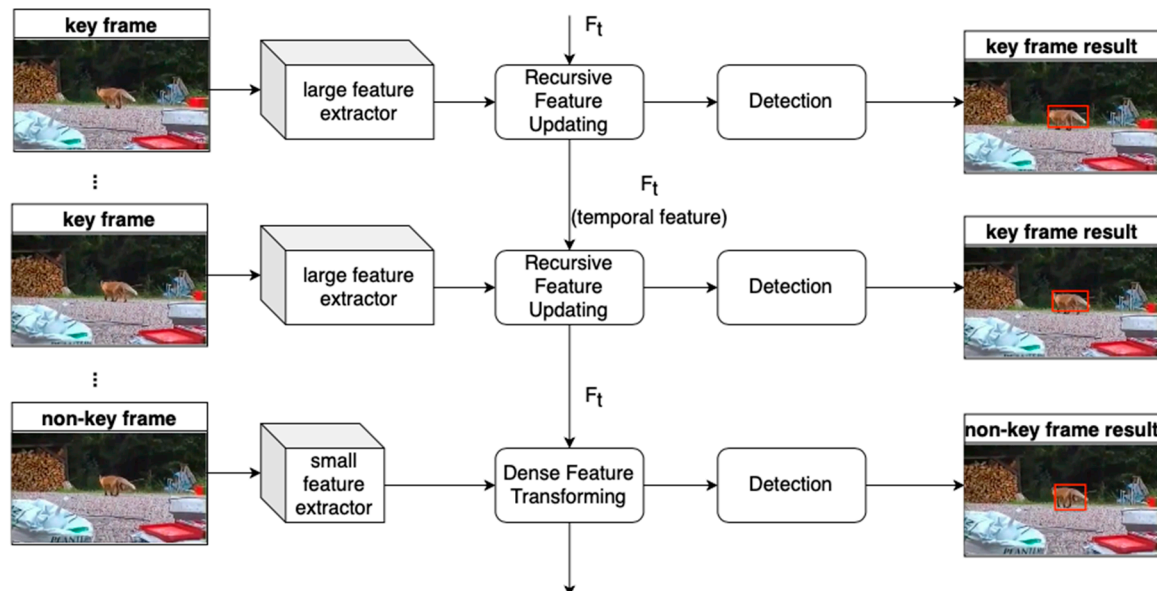


**Figure 7.** PSLA (progressive parse local attention) framework [76].

### 3.4. Tracking-Based

Inspired by the fact that tracking is an efficient way to utilize the temporal information, several methods [78,79,81] have been developed to detect objects on fixed interval frames and track them in frames in between. The improved methods in [26] and [80] detect interval frames adaptively and track the other frames.

A framework named CDT was presented in [79], combining detection and tracking for video object detection. This framework consisted of an object detector, a forward tracker and a backward tracker. Initially, objects were detected by the image object detector. Then, each detected object was tracked by the forward tracker, and undetected objects were stored by the backward tracker. In the entire process, the object detector and the tracker cooperated with each other to deal with the appearance and disappearance of objects.

Another framework named CaTDet with a high computational efficiency was presented in [78]. This framework is shown in Figure 8, which includes a tracker and a detector. CaTDet uses a tracker to predict the position of objects with a high confidence in a next frame. The processing steps of CaTDet

are: (i) Every frame is inputted to a proposal network to output potential proposals in the frame. (ii) Object position in a next frame is predicted with a high confidence using the tracker. (iii) In order to obtain the calibrated object information, the outputs of the tracker and the proposal network are combined and inputted to a refinement network.
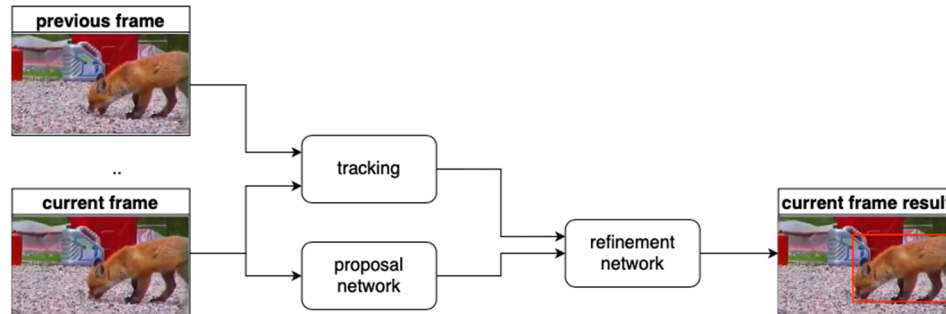


**Figure 8.** CaTDet framework [78].

More specifically, based on the observation that objects detected in one video frame would most likely appear in a next frame, a tracker was used to predict the positions on the next frame with the historical information. In case new objects appeared in a current frame, a computationally efficient proposal network similar to RPN was utilized to detect proposals. In addition, to address situations such as motion blur and occlusion, the temporal information was used by a tracker to predict future positions. The results obtained by combining the tracker and the proposal network was then refined by a refinement network. Only the regions of interest were refined by the refinement network to save computation time while maintaining accuracy.

Similar to CDT and CaTDet, recent approaches for the detection and tracking of objects in video involve rather complex multistage components. In [81], a framework using a ConvNet architecture was deployed in a simple but effective way by performing tracking and detection simultaneously. More specifically, first R-FCN [19] was employed to extract the feature maps shared between detection and tracking. Then, proposals in each frame were obtained by using RPN based on anchors [15]. RoI pooling [15] was utilized for the final detection. In particular, a regressor was introduced to extend the architecture. Position-sensitive regression maps from both frames were used together with correlation maps as the input to an RoI tracking module, in which the box relationship between the two frames was outputted. For video object detection, the framework in [81] was evaluated on the ImageNet VID dataset achieving an accuracy of 82.0% mAP.

Similarly, inspired by the observation that object tracking is more efficient than object detection, a framework (D or T) was covered in [80], see Figure 9, which includes a scheduler network to determine the operation (detecting or tracking) on a certain frame. Compared with the baseline frame skipping (detecting on fixed interval frames and tracking on intermediate frames), the scheduler network with light weights and a simple structure was found to be more effective on the ImageNet VID dataset. Moreover, the adaptive mechanism in [26] (TRACKING ASSISTED) was used to select key frames. Detection on key frames involved the utilization of an accurate detection network and detection on non-key frames was assisted by the tracking module.
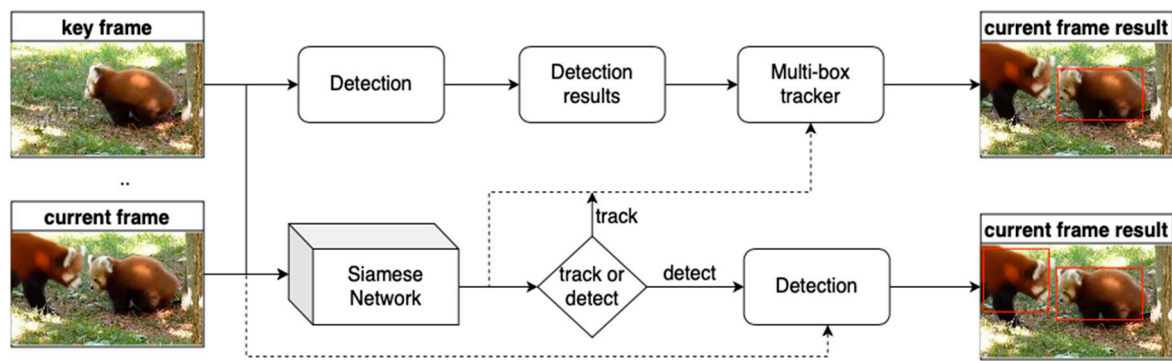
**Figure 9.** D or T framework [80].

## 3.5. Other Methods

Apart from the frameworks described above, some methods are presented that are based on a combination of multiple methods described above [24,121]. The method in [24] is based on the optical flow and tracking methods. The methods in [121] (Attentional LSTM) and [122] (TSSD) are based on the attention and LSTM methods.

In addition, these other methods appear in the literature [36,83–90]. The methods in [83] and [87] discuss ways to align and enhance feature maps. While the method in [90] studied the effect of the input image size by selecting a size to achieve a better speed-accuracy trade-off. The method in [83] named STSN (spatiotemporal sampling networks) is shown in Figure 10. This method aligns feature maps between adjacent frames. Similar to the FGFA method in [22], it relies on the idea that detection on a single frame would have difficulties dealing with noise sources such as motion blur and video defocus. Multiple frames are thus utilized for feature enhancement to achieve better performance. Unlike FGFA, which uses the optical flow method to align feature maps, deformable convolution is employed for feature alignment in [83]. First, a sharing feature extraction network is applied to extract feature maps on a current frame and adjacent frames. Then, the two feature maps are concatenated per channel and a deformable convolution is performed. The result of the deformable convolution is used as the offset for the second deformable convolution operation to align the feature maps. Furthermore, augmented feature maps are obtained by aggregating the features in the same way as FGFA. Compared with FGFA, STSN uses deformable convolution to align the features of two adjacent frames implicitly. Although it is not as intuitive as the optical flow method, it is also found to be effective. According to the experimental results reported, STSN still achieved a higher mAP than FGFA (78.9% vs. 78.8%) without relying on the optical flow information. In addition, without the assistant of the temporal post-processing, STSN obtained a better performance than the D&T baseline [81], 78.9% vs. 75.8%.
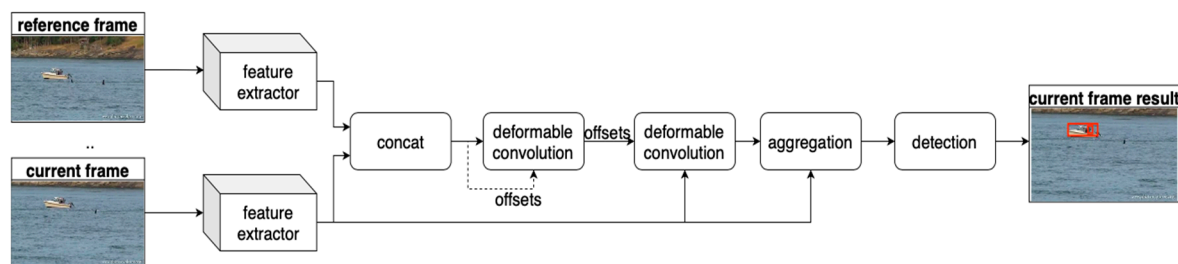


**Figure 10.** STSN (spatiotemporal sampling networks) framework [83].

Different from [83], by using the deformable convolution to propagate the temporal information, the Spatial-Temporal Memory Network (STMN) was considered in [87], which involved an RNN architecture with a Spatial-Temporal Memory module (STMM) to incorporate the long-term temporal information. The Spatial-Temporal Memory Network (STMN) operates in an end-to-end manner to model the long-term information and align the motion dynamics for video object detection. STMM

is the core module in STMN, a convolutional recurrent computation unit which fully utilizes the pretrained weights learned from static image datasets such as ImageNet [93]. This design is essential to address the practical difficulties of learning from video datasets, which largely lack the diversity of objects within the same category. STMM receives the feature maps of a current frame at time step $t$ and the spatial-temporal memory $\overrightarrow{M_{t-1}}$ with the information of all the previous frames. Then, the spatial-temporal memory $\overrightarrow{M_t}$ of the current time step is updated. In order to capture the information of both later frames and previous frames at the same time, two STMMs are used for bidirectional feature aggregation to produce the memory $M$, which is employed for both classification and bounding box regression. Therefore, the feature maps are propagated and aggregated by combining the information across multiple video frames. Evaluated on the ImageNet VID dataset, STMN has achieved the current start-of-the-art accuracy.

All the algorithms described above start from how to propagate and aggregate feature maps. In [90], video object detection was examined from another point of view. Similar to [123], the effect of input image size on the performance of video object detection was studied in [90]. Furthermore, it was found that down-sampling images can obtain better accuracy sometimes. From this point of view, a framework named AdaScale was proposed to adaptively select the input image size. AdaScale predicts the best scale or size of a next frame according to the information of a current frame. One of the reasons for the improvement is that the number of false positives is reduced. The other reason for this is that the number of true positives is increased by resizing the too-large objects to a suitable size for the detector.

In [90], the optimal scale (pixels of the shortest side) of a given image is defined with a predefined finite set of scales S (S = {600, 480, 360, 240} in [90]). Furthermore, a loss function consisting of the classification and regression loss is employed as the evaluation metric to compare the results across different scales. The regression loss for the background is expected to be zero. Hence, if the loss function is utilized directly to evaluate the results across different scales, the image scale which contains fewer foreground bounding boxes is supported. In order to deal with this issue, a new metric (the loss function, which focuses on the same number of foreground bounding boxes chosen on different scales) is employed to compare across different scales. More specifically, the number of bounding boxes involved to compute the loss is determined by the minimum number ($m$) on all the scales. For each scale, the loss of the predicted foreground bounding boxes on the image is sorted in ascending order and the first $m$ bounding boxes are chosen. The scale $m$ with the minimum loss is defined as the best scale. Inspired by R-FCN [19] working on deep features for bounding boxes regression, the channels of the deep features are expected to contain the size information. Therefore, a scale regressor using deep features is built to predict the optimal scale. Evaluated on the ImageNet VID and mini YouTube-BB datasets, Adascale achieved 1.3% and 2.7% mAP improvements with 1.6 and 1.8 times speedup compared with single-scale training and testing, respectively. Furthermore, combined with DFF [28], the speed was increased by 25% while maintaining mAP on the ImageNet VID dataset.

## 4. Comparison of Video Object Detection Methods

The great majority of video object detection approaches use the ImageNet VID dataset [5] for performance evaluation. In this section, the timeline of video object detection methods in recent years is shown in Figure 11 together with a group listing of the methods in Figure 12. Then, a comparison is provided between the methods covered in the previous section. The comparison is presented in Tables 4 and 5, which correspond to with and without post-processing, respectively. The methods in Figure 11 belong to different groups but the same time, whereas the methods in Figure 12 belong to different times but the same groups. As can be seen from Figures 11 and 12, the methods based on optical flow were proposed earlier. During the same period, video object detection methods were assisted by tracking due to the effectiveness of tracking in utilizing the temporal–spatial information. The optical flow-based methods needed a large number of parameters and they were only suitable for small motions. In recent years, the methods based on attention have achieved much success, such as

MEGA [25]. Using LSTM for feature propagation and aggregation is becoming a hot research topic and many new methods are being proposed, such as STSN [83] using deformable convolution to align the feature maps. The latest research is mostly based on attention, LSTM or a combination of methods such as Flow&LSTM [72].
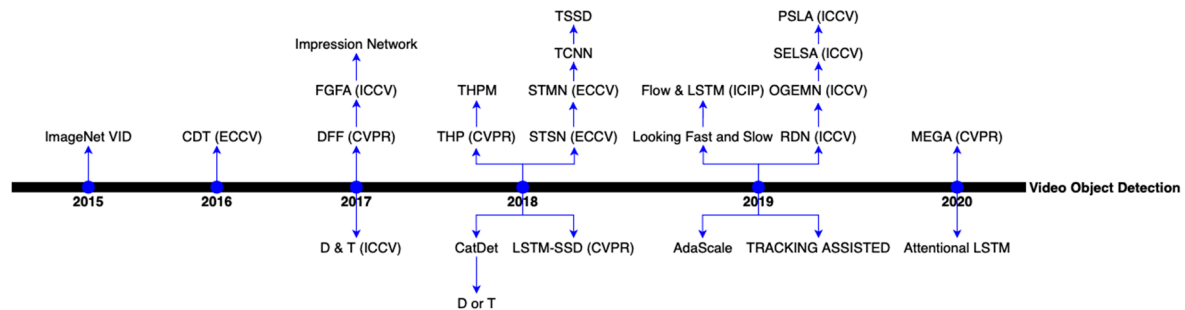


**Figure 11.** Timeline of video object detection methods.
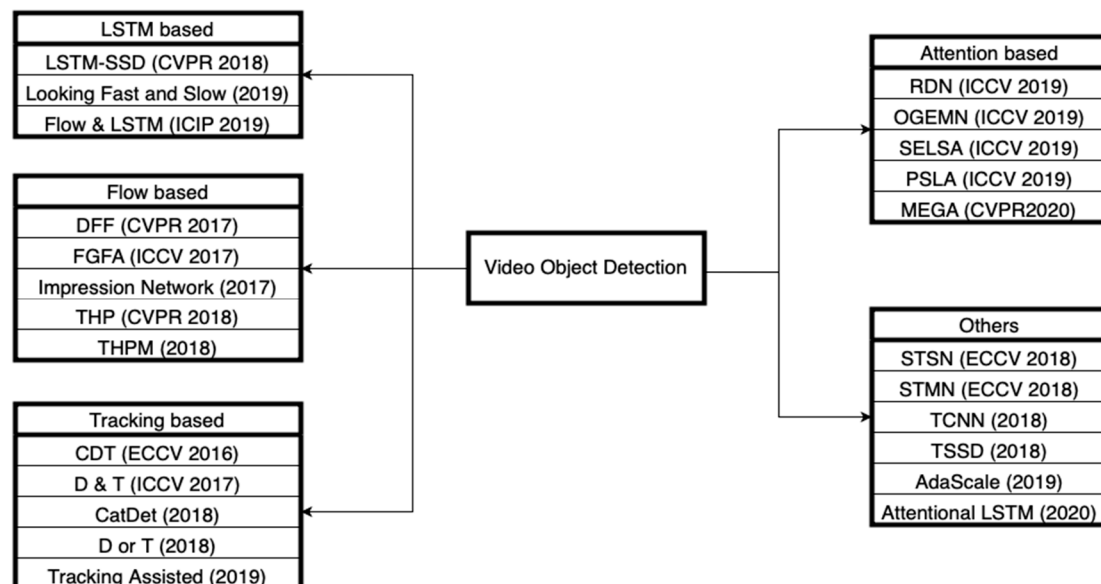


**Figure 12.** Video object detection methods sorted in different groups.

Table 4 provides the outcomes without post processing. In this table, the methods are divided into different groups according to the way temporal and spatial information are utilized. Flow-guided groups propagate and align the feature maps according to the flow field obtained by optical flow. Both accuracy and speed of various frameworks are reported in this table. For example, DFF provides high computational efficiency and achieves a runtime of 20.25 fps using a Titan K40 GPU. FGFA achieves a high accuracy producing 76.3% mAP with 1.36 fps. Obviously, DFF is faster than FGFA. Flow-guided methods are intuitive and well understood to propagate features. Optical flow is deemed suitable for small movement estimation. In addition, since optical flow reflects pixel level displacement, it has difficulties when it is applied to high-level feature maps. One pixel movement on feature maps may correspond to 10 to 20 pixels movement.

**Table 4.** Comparison among the video object detection methods without post processing; note that
the runtime is based on the NVIDIA GPU used in the references: K means K40, XP means Titan XP,
X means Titan X, V means Titan V, 1060 means GeForce GTX 1060, 1080 Ti means GeForce GTX 1080 Ti,
2080 Ti means GeForce GTX 2080 Ti.

| Type | Framework | Backbone | mAP (%) | Runtime (fps) |
|---|---|---|---|---|
| Single frame | R-RCN [19] | ResNet-101 | 73.9 | 4.05 K |
| | | | 70.3 | 12 XP |
| | Impression network [68] | Modified ResNet-101 | 75.5 | 20 1060 |
| Flow based | FGFA [22] | ResNet-101 | 76.3 | 1.36 K |
| | DFF [28] | ResNet-101 | 73.1 | 20.25 K |
| | THP [69] | ResNet-101 + DCN | 78.6 | 13.0X/8.6K |
| | THPM [27] | Mobilenet | 60.2 | 25.6 HuaiWei Mate8 |
| LSTM based | Looking fast and slow [70] | Interleaved | 61.4 | 23.5 Pixel 3 phone |
| | LSTM-SSD [71] | MobilenetV2-SSDLite | 53.5 | − |
| | Flow&LSTM [72] | ResNet-101 | 75.5 | − |
| | OGEMN [74] | ResNet-101 | 79.3 | 8.9 (1080Ti) |
| | | ResNet-101 + DCN | 80.0 | − |
| Attention based | PSLA [76] | ResNet-101 | 77.1 | 30.8V\18.73X |
| | | ResNet-101 + DCN | 80.0 | 26.0V\13.34X |
| | SELSA [77] | ResNet-101 | 80.25 | − |
| | | ResNeXt-101 | 83.11 | |
| | RDN [75] | ResNet-101 | 81.8 | 10.6 V100 |
| | | ResNeXt-101 | 83.2 | − |
| | MEGA [25] | ResNet-101 | 82.9 | 8.73 2080Ti |
| | | ResNeXt-101 | 84.1 | − |
| Tracking based | D&T loss [81] | ResNet-101 | 75.8 | 7.8X |
| | Track assisted [26] | ResNet-101 | 70.0 | 30XP |
| Others | TCNN [24] | GoogLeNet | 73.8 | − |
| | STSN [83] | ResNet-101 + DCN | 78.9 | − |

**Table 5.** Comparison among the video object detection methods with post processing.

| Type | Framework | Backbone | mAP (%) | Runtime (fps) |
|---|---|---|---|---|
| Flow-based | FGFA [22] | ResNet-101 | 78.4 | − |
| | | Inception-ResNet | 80.1 | |
| | Looking fast and slow [70] | Interleaved + Quantization + Async | 59.3 | 72.3 Pixel 3 phone |
| Lstm-based | MobilenetV2-SSDLite + LSTM ($\alpha = 1.4$) [71] | MobilenetV2-SSDLite | 64.1 | 4.1 Pixel 3 phone |
| | MobilenetV2-SSDLite + LSTM ($\alpha = 1.0$) [71] | MobilenetV2-SSDLite | 59.1 | − |
| | MobilenetV2-SSDLite + LSTM ($\alpha = 0.5$) [71] | MobilenetV2-SSDLite | 50.3 | − |
| | MobilenetV2-SSDLite + LSTM ($\alpha = 0.35$) [71] | MobilenetV2-SSDLite | 45.1 | 14.6 Pixel 3 phone |
| | OGEMN [74] | ResNet-101 | 80.8 | − |
| | | ResNet-101 + DCN | 81.6 | |
| Attention-based | PSLA [76] | ResNet-101 | 78.6 | 5.7X |
| | | ResNet-101 + DCN | 81.4 | 6.31V\5.13X |
| | SELSA [77] | ResNet-101 | 80.54 | − |
| | RDN [75] | ResNet-101 | 83.8 | − |
| | | ResNeXt-101 | 84.7 | |
| | MEGA [25] | ResNet-101 | 84.5 | − |
| | | ResNeXt-101 | 85.4 | |
| Tracking-based | D&T ($\tau = 10$) [81] | ResNet-101 | 78.6 | − |
| | D&T ($\tau = 1$) [81] | ResNet-101 | 79.8 | 5X |
| | D&T [81] | Inception V4 | 82.0 | − |
| Others | STSN [83] | ResNet-101 + DCN | 80.4 | − |
| | STMN [87] | ResNet-101 | 80.5 | − |

Inspired by the LSTM-based solutions in natural language processing, LSTM methods are used to incorporate the sequence information. In the LSTM group, Flow&LSTM [72] achieved the highest accuracy of 75.5%. Looking Fast and Slow [70] generated high speed but with low accuracy. LSTM captures the long-term information with a simple implementation. Since the sigmoid activation of the input and forget gates are rarely completely saturated, a slow state decay and thus loss of long-term dependence is resulted. In other words, it is difficult to retain the complete previous state in the update.

Attention-based methods also show the ability to perform video object detection effectively. In the attention-related group, MEGA [25] with ResNeXt-101 as backbone achieved the highest accuracy of 84.1% mAP. As described, it achieved a very high accuracy with a relatively fast speed. Attention-based methods aggregate the features within proposals that are generated. This decreases the computation

time. Because of only using the features within the proposals, the performance relies on the effect of RPN to a certain extent. Here, it is rather difficult to utilize more comprehensive information.

In the tracking-based group, the methods are assisted by tracking. D&T loss [81] achieved 75.8% mAP. Tracking is an efficient method to employ the temporal information with a detector assisted by a tracker. However, it cannot solve the problems created by motion blur and video defocus directly. As the detection performance relies on the tracking performance, the detector part suffers from tracking errors. There are also other standalone methods including TCNN [24], STSN [83] and STMN [87].

In order to further improve the performance in terms of detection accuracy, post-processing can be added to the above methods. The results with post-processing are shown in Table 5. One can easily see that with post-processing, the accuracy is noticeably improved. For example, the accuracy of MEGA is improved from 84.1% to 85.4% mAP.

## 5. Future Trends

Challenges still remain for further improving the accuracy and speed of the video object detection methods. This section presents the major challenges and possible future trends as related to video object detection.

At present, there is a lack of a comprehensive benchmark dataset containing the labels of each frame. The most widely used dataset, that is ImageNet VID, does not include complex real-world conditions as compared to the static image dataset COCO. The number of objects in each frame in the ImageNet VID dataset is limited, which is not the case under real-world conditions. In addition, in many real-world applications, videos include a large field of view and in some cases high resolution images. Lack of a well-annotated dataset representing actual or real-world conditions remains a challenge for the purpose of advancing video object detection. Hence, the establishment of a comprehensive benchmark dataset is considered a future trend of importance.

Up to now, the most widely used evaluation metric in video object detection is mAP, which is derived from static image object detection. This metric does not fully reflect the temporal characteristics in video object detection. Although Average Delay (AD) is proposed to reflect the temporal characteristics, it is still not a fully developed metric. For example, the stability of detection in video is not reflected by it. Therefore, novel evaluation metrics to reflect detection stability which are more suitable for video object detection are considered to be another future trend of importance.

Most of the methods covered in this review paper only utilize the local temporal information or global information separately. There are only a few methods, such as MEGA, which have used the local and global temporal information at the same time and achieved a benchmark mAP of 85.4%. As demonstrated by MEGA, it is worth developing future frameworks which utilize both the local and global temporal information. Furthermore, for most of the existing video object detection algorithms, the number of frames used is too small to fully utilize the video information. Hence, as yet another future trend, it is of importance to develop methods that utilize the long-term video information.

As can be observed from Tables 4 and 5, the attention-based frameworks achieved a relatively high accuracy. However, such methods pose difficulties for real-time applications demanding very powerful GPUs. Although the Looking Fast and Slow method [70] achieved 72.3 fps on Pixel 3 phones, the accuracy is only 59.3%, which poses challenges for actual deployment. Indeed, the trade-off between accuracy and speed needs to be further investigated. Real-time performance is important for practical applications such as autonomous driving and video surveillance. It is significant to pay more attention to the methods to make a light model, while ensuring that the accuracy will not drop too much. Some light network structure design methods like Depthwise Separable Convolution [115] and channel shuffle [124] used in the classification application can be used for reference in video object detection. In addition, model compression methods like [125] can be considered as well.

## 6. Conclusions

In recent years, after the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) announced the video object detection task in 2015, many deep learning-based video object detection solutions have been developed. This paper has provided a review of the video object detection methods that have been developed so far. This review has covered the available datasets, evaluation metrics and an overview of different categories of deep learning-based methods for video object detection. A categorization of the video object detection methods has been made according to the way temporal and spatial information are used. These categories include flow-based, LSTM-based, attention-based and tracking-based methods, as well as others. The performance of various detectors with or without post-processing is summarized in Tables 4 and 5 in terms of both detection accuracy and computation speed. Several trends of importance in video object detection have also been stated for possible future works.

## References

1. Bateni, S.; Wang, Z.; Zhu, Y.; Hu, Y.; Liu, C. Co-Optimizing Performance and Memory Footprint Via Integrated CPU/GPU Memory Management, an Implementation on Autonomous Driving Platform. In Proceedings of the 2020 IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS), Sydney, Australia, 21–24 April 2020.
2. Lu, J.; Tang, S.; Wang, J.; Zhu, H.; Wang, Y. A Review on Object Detection Based on Deep Convolutional Neural Networks for Autonomous Driving. In Proceedings of the 2019 Chinese Control and Decision Conference (CCDC), Nanchang, China, 3–5 June 2019.
3. Wei, H.; Laszewski, M.; Kehtarnavaz, N. Deep Learning-Based Person Detection and Classification for Far Field Video Surveillance. In Proceedings of the 2018 IEEE 13th Dallas Circuits and Systems Conference, Dallas, TX, USA, 12 November 2018.
4. Guillermo, M.; Tobias, R.R.; De Jesus, L.C.; Billones, R.K.; Sybingco, E.; Dadios, E.P.; Fillone, A. Detection and Classification of Public Security Threats in the Philippines Using Neural Networks. In Proceedings of the 2020 IEEE 2nd Global Conference on Life Sciences and Technologies (LifeTech), Kyoto, Japan, 10–12 March 2020; pp. 1–4.
5. Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; et al. ImageNet Large Scale Visual Recognition Challenge. *Int. J. Comput. Vis.* **2015**, *115*, 211–252. [CrossRef]
6. Shen, Z.; Liu, Z.; Li, J.; Jiang, Y.-G.; Chen, Y.; Xue, X. DSOD: Learning Deeply Supervised Object Detectors from Scratch. In Proceedings of the 2017 IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 1937–1945.
7. Tian, Z.; Shen, C.; Chen, H.; He, T. FCOS: Fully Convolutional One-Stage Object Detection. In Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, Korea, 27 October–2 November 2019.
8. Zhao, Q.; Sheng, T.; Wang, Y.; Tang, Z.; Chen, Y.; Cai, L.; Ling, H. M2Det: A Single-Shot Object Detector Based on Multi-Level Feature Pyramid Network. In Proceedings of the AAAI Conference on Artificial Intelligence, Honolulu, HI, USA, 27 January–1 February 2019; pp. 9259–9266.
9. Liu, W.; Anguelov, D.; Erhan, D.; Szegedy, C.; Reed, S.; Fu, C.-Y.; Berg, A.C. SSD: Single Shot MultiBox Detector. In *Computer Vision—Eccv 2016*; Part I; Leibe, B., Matas, J., Sebe, N., Welling, M., Eds.; Springer: Cham, Switzerland, 2016; pp. 21–37.
10. Redmon, J.; Farhadi, A. YOLO9000: Better, Faster, Stronger. In Proceedings of the 30th IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2016; pp. 6517–6525.
11. Redmon, J.; Farhadi, A. YOLOv3: An. Incremental Improvement. *arXiv* **2018**, arXiv:1804.02767.

12. Redmon, J.; Divvala, S.; Girshick, R.; Farhadi, A. You Only Look Once: Unified, Real-Time Object Detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 779–788.
13. He, K.; Gkioxari, G.; Dollar, P.; Girshick, R. Mask R-CNN. In Proceedings of the IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 2980–2988.
14. Girshick, R. Fast R-CNN. In Proceedings of the 2015 IEEE International Conference on Computer Vision, Santiago, Chile, 7–13 December 2015; pp. 1440–1448.
15. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster R-CNN: Towards Real-Time Object Detection with Region. Proposal Networks. *IEEE Trans. Pattern Anal. Mach. Intell.* **2017**, *39*, 1137–1149. [CrossRef]
16. Girshick, R.; Donahue, J.; Darrell, T.; Malik, J. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition, Columbus, OH, USA, 23–28 June 2014; pp. 580–587.
17. Cai, Z.; Vasconcelos, N. Cascade R-CNN: Delving into High Quality Object Detection. In Proceedings of the 2018 IEEE/Cvf Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–22 June 2018; pp. 6154–6162.
18. He, K.; Zhang, X.; Ren, S.; Sun, J. Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition. In *Computer Vision—Eccv 2014*; Part III; Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T., Eds.; IEEE: Zurich, Switzerland, 2014; pp. 346–361.
19. Dai, J.; Li, Y.; He, K.; Sun, J. R-FCN: Object Detection via Region.-based Fully Convolutional Networks. In *Advances in Neural Information Processing Systems 29*; Lee, D.D., Sugiyama, M., Luxburg, U.V., Guyon, I., Garnett, R., Eds.; Curran Associates, Inc.: Red Hook, NY, USA, 2016.
20. Shrivastava, A.; Gupta, A.; Girshick, R. Training Region—Based Object Detectors with Online Hard Example Mining. In Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 761–769.
21. Wei, H.; Kehtarnavaz, N. Semi-Supervised Faster RCNN-Based Person Detection and Load Classification for Far Field Video Surveillance. *Mach. Learn. Knowl. Extr.* **2019**, *1*, 756–767. [CrossRef]
22. Zhu, X.; Wang, Y.; Dai, J.; Yuan, L.; Wei, Y. Flow-Guided Feature Aggregation for Video Object Detection. In Proceedings of the 2017 IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 408–417.
23. Zhang, R.; Miao, Z.; Zhang, Q.; Hao, S.; Wang, S. Video Object Detection by Aggregating Features across Adjacent Frames. In Proceedings of the 2019 3rd International Conference on Machine Vision and Information Technology, Guangzhou, China, 22–24 February 2019. [CrossRef]
24. Kang, K.; Li, H.; Yan, J.; Zeng, X.; Yang, B.; Xiao, T.; Zhang, C.; Wang, Z.; Wang, R.; Wang, X.; et al. T-CNN: Tubelets With Convolutional Neural Networks for Object Detection from Videos. *IEEE Trans. Circuits Syst. Video Technol.* **2018**, *28*, 2896–2907. [CrossRef]
25. Chen, Y.; Cao, Y.; Hu, H.; Wang, L. Memory Enhanced Global-Local Aggregation for Video Object Detection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Seattle, WA, USA, 16–18 June 2020.
26. Yang, W.; Liu, B.; Li, W.; Yu, N. Tracking Assisted Faster Video Object Detection. In Proceedings of the 2019 IEEE International Conference on Multimedia and Expo, Shanghai, China, 8–12 July 2019; pp. 1750–1755.
27. Zhu, X.; Dai, J.; Zhu, X.; Wie, Y.; Yuan, L. Towards High Performance Video Object Detection for Mobiles. *arXiv* **2018**, arXiv:1804.05830.
28. Zhu, X.; Xiong, Y.; Dai, J.; Yuan, L.; Wei, Y. Deep Feature Flow for Video Recognition. In Proceedings of the 30th IEEE Conference on Computer Vision and Pattern Recognition, Venice, Italy, 22–29 October 2017; pp. 4141–4150.
29. Horn, B.K.; Schunck, B.G. Determining Optical-Flow. *Artif. Intell.* **1981**, *17*, 185–203. [CrossRef]
30. Nguyen, H.T.; Worring, M.; Dev, A. Detection of moving objects in video using a robust motion similarity measure. *IEEE Trans. Image Process.* **2000**, *9*, 137–141. [CrossRef] [PubMed]
31. Carminati, L.; Benois-Pineau, J. Gaussian mixture classification for moving object detection in video surveillance environment. In Proceedings of the 2005 International Conference on Image Processing, Genova, Italy, 11–14 September 2005; pp. 3361–3364.

32. Jayabalan, E.; Krishnan, A. Object Detection and Tracking in Videos Using Snake and Optical Flow Approach. In *Computer Networks and Information Technologies*; Das, V.V., Stephen, J., Chaba, Y., Eds.; Springer: Berlin/Heidelberg, Germany, 2011; p. 299.

33. Jayabalan, E.; Krishnan, A. Detection and Tracking of Moving Object in Compressed Videos. In *Computer Networks and Information Technologies*; Das, V.V., Stephen, J., Chaba, Y., Eds.; Springer: Berlin/Heidelberg, Germany, 2011; pp. 39–43.

34. Ghosh, A.; Subudhi, B.N.; Ghosh, S. Object Detection from Videos Captured by Moving Camera by Fuzzy Edge Incorporated Markov Random Field and Local Histogram Matching. *IEEE Trans. Circuits Syst. Video Technol.* **2012**, *22*, 1127–1135. [CrossRef]

35. Guo, C.; Gao, H. Adaptive graph-cuts algorithm based on higher-order MRF for video moving object detection. *Electron. Lett.* **2012**, *48*, 371–373. [CrossRef]

36. Guo, C.; Liu, D.; Guo, Y.; Sun, Y. An adaptive graph cut algorithm for video moving objects detection. *Multimed. Tools Appl.* **2014**, *72*, 2633–2652. [CrossRef]

37. Yadav, D.K.; Singh, K. A combined approach of Kullback-Leibler divergence and background subtraction for moving object detection in thermal video. *Infrared Phys. Technol.* **2016**, *76*, 21–31. [CrossRef]

38. Oreifej, O.; Li, X.; Shah, M. Simultaneous Video Stabilization and Moving Object Detection in Turbulence. *IEEE Trans. Pattern Anal. Mach. Intell.* **2013**, *35*, 450–462. [CrossRef]

39. Nadimi, S.; Bhanu, B. Physical models for moving shadow and object detection in video. *IEEE Trans. Pattern Anal. Mach. Intell.* **2004**, *26*, 1079–1087. [CrossRef]

40. Utsumi, O.; Miura, K.; Ide, I.; Sakai, S.; Tanaka, H. An object detection method for describing soccer games from video. In Proceedings of the IEEE International Conference on Multimedia and Expo, Lausanne, Switzerland, 26–29 August 2002; pp. 45–48.

41. Hossain, M.J.; Dewan MA, A.; Chae, O. Moving object detection for real time video surveillance: An. Edge based approach. *IEICE Trans. Commun.* **2007**, *90*, 3654–3664. [CrossRef]

42. Chiranjeevi, P.; Sengupta, S. Robust detection of moving objects in video sequences through rough set theory framework. *Image Vis. Comput.* **2012**, *30*, 829–842. [CrossRef]

43. Abd Razak, H.; Abd Almisreb, A.; Saleh, M.A.; Tahir, N.M. Anomalous Behaviour Detection using Transfer Learning Algorithm of Series and DAG Network. In Proceedings of the 2019 IEEE 9th International Conference on System Engineering and Technology, Shah Alam, Malaysia, 7 October 2019; pp. 505–509.

44. Azarang, A.; Manoochehri, H.E.; Kehtarnavaz, N. Convolutional Autoencoder-Based Multispectral Image Fusion. *IEEE Access* **2019**, *7*, 35673–35683. [CrossRef]

45. Majumder, S.; Elloumi, Y.; Akil, M.; Kachouri, R.; Kehtarnavaz, N. A deep learning-based smartphone app for real-time detection of five stages of diabetic retinopathy. In Proceedings of the Real-Time Image Processing and Deep Learning 2020, Online Only, CA, USA, 27 April–8 May 2020.

46. Wang, Z.; Wang, Y.; Lin, Y.; Delord, E.; Latifur, K. Few-Sample and Adversarial Representation Learning for Continual Stream Mining. In Proceedings of the WWW '20: The Web Conference 2020, Taipei, Taiwan, 20–24 April 2020.

47. Maor, G.; Zeng, X.; Wang, Z.; Hu, Y. An FPGA Implementation of Stochastic Computing-based LSTM. In Proceedings of the 2019 IEEE 37th International Conference on Computer Design, Abu Dhabi, UAE, 17–20 November 2019; pp. 38–46.

48. Chu, X. Human Pose Estimation and Immediacy Prediction with Deep Learning. Ph.D. Thesis, The Chinese University of Hong Kong, Hong Kong, China, August 2017.

49. Wang, Z.; Tao, H.; Kong, Z.; Chandra, S.; Khan, L. Metric Learning based Framework for Streaming Classification with Concept Evolution. In Proceedings of the 2019 International Joint Conference on Neural Networks, Budapest, Hungary, 14–19 July 2019.

50. Li, H.; Meng, L.; Zhang, J.; Tan, Y.; Ren, Y.; Zhang, H. Multiple Description Coding Based on Convolutional Auto-Encoder. *IEEE Access* **2019**, *7*, 26013–26021. [CrossRef]

51. Zheng, S.; Liu, G.; Suo, H.; Lei, Y. Autoencoder-Based Semi-Supervised Curriculum Learning for Out-of-Domain Speaker Verification. In Proceedings of the INTERSPEECH 2019, Graz, Austria, 15–19 September 2019; pp. 4360–4364. [CrossRef]

52. Wei, H.; Kehtarnavaz, N. Determining Number of Speakers from Single Microphone Speech Signals by Multi-Label. Convolutional Neural Network. In Proceedings of the IECON 2018—44th Annual Conference of the IEEE Industrial Electronics Society, Washington, DC, USA, 21–23 October 2018; pp. 2706–2710.

53. Zhao, Y.; Wang, D.; Merks, I.; Zhang, T. Dnn-Based Enhancement of Noisy and Reverberant Speech. In Proceedings of the 2016 IEEE International Conference on Acoustics, Speech and Signal, Shanghai, China, 20–25 March 2016; pp. 6525–6529.

54. Tao, F.; Liu, G.; Zhao, Q. An Ensemble Framework of Voice-Based Emotion Recognition System for Films and Tv Programs. In Proceedings of the 2018 IEEE International Conference on Acoustics, Speech and Signal Processing, Calgary, AB, Canada, 15–20 April 2018; pp. 6209–6213.

55. Zhao, Y.; Xu, B.; Giri, R.; Zhang, T. Perceptually Guided Speech Enhancement Using Deep Neural Networks. In Proceedings of the 2018 IEEE International Conference on Acoustics, Speech and Signal Processing, Calgary, AB, Canada, 15–20 April 2018; pp. 5074–5078.

56. Tao, F.; Busso, C. Aligning Audiovisual Features for Audiovisual Speech Recognition. In Proceedings of the IEEE International Conference on Multimedia and Expo, San Diego, CA, USA, 23–27 July 2018.

57. Wei, H.; Chopada, P.; Kehtarnavaz, N. C-MHAD: Continuous Multimodal Human Action Dataset of Simultaneous Video and Inertial Sensing. *Sensors* **2020**, *20*, 2905. [CrossRef]

58. Brena, R.F.; Aguileta, A.A.; Trejo, L.A.; Molino-Minero-Re, E.; Mayora, O. Choosing the Best Sensor Fusion Method: A Machine-Learning Approach. *Sensors* **2020**, *20*, 2350. [CrossRef]

59. Tao, F.; Busso, C. End-to-End Audiovisual Speech Recognition System with Multitask Learning. *IEEE Trans. Multimed.* **2020**. [CrossRef]

60. Wei, H.; Kehtarnavaz, N. Simultaneous Utilization of Inertial and Video Sensing for Action Detection and Recognition in Continuous Action Streams. *IEEE Sens. J.* **2020**, *20*, 6055–6063. [CrossRef]

61. Chen, C.; Jafari, R.; Kehtarnavaz, N. A survey of depth and inertial sensor fusion for human action recognition. *Multimed. Tools Appl.* **2017**, *76*, 4405–4425. [CrossRef]

62. Li, M.; Sun, L.; Huo, Q. Dff-Den: Deep Feature Flow with Detail Enhancement Network for Hand Segmentation in Depth Video. In Proceedings of the 2018 25th IEEE International Conference on Image Processing, Athens, Greece, 7–10 October 2018; pp. 1548–1552.

63. Li, M.; Sun, L.; Huo, Q. Flow-guided feature propagation with occlusion aware detail enhancement for hand segmentation in egocentric videos. *Comput. Vis. Image Underst.* **2019**, *187*. [CrossRef]

64. Li, H.; Yang, W.; Liao, Q. Temporal Feature Enhancing Network for Human Pose Estimation in Videos. In Proceedings of the 2019 IEEE International Conference on Image Processing, Taipei, Taiwan, 22–25 September 2019; pp. 579–583.

65. Zhou, Q.; Liang, X.; Gong, K.; Lin, L. Adaptive Temporal Encoding Network for Video Instance-level Human Parsing. In Proceedings of the 2018 ACM Multimedia Conference, Seoul, Korea, 22–26 October 2018; pp. 1527–1535.

66. Pi, Z.; Qin, H.; Gao, C.; Sang, N. Jointly detecting and multiple people tracking by semantic and scene information. *Neurocomputing* **2020**, *412*, 244–251. [CrossRef]

67. Wang, S.; Zhou, Y.; Yan, J.; Deng, Z. Fully Motion-Aware Network for Video Object Detection. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8–14 September 2018; Springer: Cham, Switzerland, 2018.

68. Hetang, C.; Qin, H.; Liu, S.; Yan, J. Impression Network for Video Object Detection. *arXiv* **2017**, arXiv:1712.05896.

69. Zhu, X.; Dai, J.; Yuan, L.; Wei, Y. Towards High Performance Video Object Detection. In Proceedings of the 2018 IEEE/Cvf Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–23 June 2018; pp. 7210–7218.

70. Liu, M.; Zhu, M.; White, M.; Li, Y.; Kalenichenko, D. Looking Fast and Slow: Memory-Guided Mobile Video Object Detection. *arXiv* **2019**, arXiv:1903.10172.

71. Liu, M.; Zhu, M. Mobile Video Object Detection with Temporally-Aware Feature Maps. In Proceedings of the 2018 IEEE/Cvf Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–23 June 2018; pp. 5686–5695.

72. Zhang, C.; Kim, J. Modeling Long—And Short-Term Temporal Context for Video Object Detection. In Proceedings of the 2019 IEEE International Conference on Image Processing, Taipei, Taiwan, 22–25 September 2019; pp. 71–75.

73. Lu, Y.; Lu, C.; Tang, C.-K. Online Video Object Detection using Association LSTM. In Proceedings of the 2017 IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 2363–2371.

74. Deng, H.; Hua, Y.; Song, T.; Zhang, Z.; Xue, Z.; Ma, R.; Guan, H. Object Guided External Memory Network for Video Object Detection. In Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, Korea, 27 October–2 November 2019; pp. 6677–6686.

75. Deng, J.; Pan, Y.; Yao, T.; Zhou, W.; Li, H.; Mei, T. Relation Distillation Networks for Video Object Detection. In Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, Korea, 27 October–2 November 2019; pp. 7022–7031.

76. Guo, C.; Fan, B.; Gu, J.; Zhang, Q.; Xiang, S.; Prinet, V.; Pan, C. Progressive Sparse Local Attention for Video object detection. In Proceedings of the IEEE International Conference on Computer Vision, Seoul, Korea, 27 October–2 November 2019.

77. Wu, H.; Chen, Y.; Wang, N.; Zhang, Z. Sequence Level Semantics Aggregation for Video Object Detection. In Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Gangnam-gu, Seoul, Korea, 27 October–2 November 2019.

78. Mao, H.; Kong, T.; Dally, W.J. CaTDet: Cascaded Tracked Detector for Efficient Object Detection from Video. *arXiv* **2018**, arXiv:1810.00434.

79. Kim, H.U.; Kim, C.S. CDT: Cooperative Detection and Tracking for Tracing Multiple Objects in Video Sequences. In *Computer Vision—Eccv 2016*; Part VI; Leibe, B., Matas, J., Sebe, N., Welling, M., Eds.; Springer: Cham, Switzerland, 2016; pp. 851–867.

80. Luo, H.; Xie, W.; Wang, X.; Zeng, W. Detect or Track: Towards Cost-Effective Video Object Detection/Tracking. In Proceedings of the AAAI Conference on Artificial Intelligence, New Orleans, LA, USA, 2–7 February 2018.

81. Feichtenhofer, C.; Pinz, A.; Zisserman, A. Detect to Track and Track to Detect. In Proceedings of the 2017 IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 3057–3065.

82. Sharma, V.K.; Acharya, B.; Mahapatra, K.K. Online Training of Discriminative Parameter for Object Tracking-by-Detection in a Video. In *Soft Computing in Data Analytics*; Nayak, J., Abraham, A., Krishna, B., Chandra Sekhar, G., Das, A., Eds.; Springer: Singapore, 2019; pp. 215–223.

83. Bertasius, G.; Torresani, L.; Shi, J. Object Detection in Video with Spatiotemporal Sampling Networks. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8–14 September 2018.

84. Chen, K.; Chen, K.; Wang, J.; Yang, S.; Zhang, X.; Xiong, Y.; Loy, C.C.; Lin, D. Optimizing Video Object Detection via a Scale-Time Lattice. In Proceedings of the 2018 IEEE/Cvf Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–22 June 2018; pp. 7814–7823.

85. Wang, T.; Xiong, J.; Xu, X.; Shi, Y. SCNN: A General Distribution Based Statistical Convolutional Neural Network with Application to Video Object Detection. In Proceedings of the AAAI Conference on Artificial Intelligence, Honolulu, HI, USA, 27 January–1 February 2019; pp. 5321–5328.

86. Du, Y.; Yuan, C.; Hu, W.; Maybank, S. Spatio-temporal self-organizing map deep network for dynamic object detection from videos. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017.

87. Xiao, F.; Jae Lee, Y. Video Object Detection with an Aligned Spatial-Temporal Memory. *arXiv* **2017**, arXiv:1712.06317.

88. Jiang, Z.; Gao, P.; Guo, C.; Zhang, Q.; Xiang, S.; Pan, C. Video Object Detection with Locally-Weighted Deformable Neighbors. In Proceedings of the AAAI Conference on Artificial Intelligence, Honolulu, HI, USA, 27 January–1 February 2019; pp. 8529–8536.

89. Zhu, H.; Yan, X.; Tang, H.; Chang, Y.; Li, B.; Yuan, X. Moving Object Detection with Deep CNNs. *IEEE Access* **2020**, *8*, 29729–29741. [CrossRef]

90. Chin, T.W.; Ding, R.; Marculescu, D. AdaScale: Towards Real-time Video Object Detection Using Adaptive Scaling. *arXiv* **2019**, arXiv:1902.02910.

91. Rybak, Ł.; Dudczyk, J. A Geometrical Divide of Data Particle in Gravitational Classification of Moons and Circles Data Sets. *Entropy* **2020**, *22*, 1088. [CrossRef]

92. Lin, T.-Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollar, P.; Zitnick, C.L. Microsoft COCO: Common Objects in Context. In *Computer Vision—Eccv 2014*; Part V; Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T., Eds.; Springer: Cham, Switzerland, 2014; pp. 740–755.

93. Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; Li, F.-F. ImageNet: A Large-Scale Hierarchical Image Database. In Proceedings of the Cvpr: 2009 IEEE Conference on Computer Vision and Pattern Recognition, Miami, FL, USA, 20–25 June 2009; pp. 248–255.

94.  Real, E.; Shlens, J.; Mazzocchi, S.; Pan, X.; Vanhoucke, V. YouTube-BoundingBoxes: A Large High—Precision Human-Annotated Data Set for Object Detection in Video. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 7464–7473.

95.  Damen, D.; Doughty, H.; Farinella, G.; Fidler, S.; Furnari, A.; Kazakos, E.; Wray, M. The Epic-Kitchens Dataset: Collection, Challenges and Baselines. *IEEE Trans. Pattern Anal. Mach. Intell.* **2020**, *1*, 1. [CrossRef]

96.  Perazzi, F.; Pont-Tuset, J.; McWilliams, B.; Van Gool, L.; Gross, M.; Sorkine-Hornung, A. A Benchmark Dataset and Evaluation Methodology for Video Object Segmentation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 724–732.

97.  Wang, Y.; Jodoin, P.-M.; Porikli, F.; Konrad, J.; Benezeth, Y.; Ishwar, P. CDnet 2014: An Expanded Change Detection Benchmark Dataset. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, Columbus, OH, USA, 23–28 June 2014; p. 393.

98.  Kristan, M.; Matas, J.; Leonardis, A.; Felsberg, M.; Cehovin, L.; Fernandez, G.; Vojir, T.; Hager, G.; Nebehay, G.; Pflugfelder, R.; et al. The Visual Object Tracking VOT2015 challenge results. In Proceedings of the IEEE International Conference on Computer Vision Workshop, Santiago, Chile, 7–13 December 2015; pp. 564–586.

99.  Leal-Taixé, L.; Milan, A.; Reid, I.; Roth, S.; Schindler, K. MOTChallenge 2015: Towards a Benchmark for Multi-Target. Tracking. *arXiv* **2015**, arXiv:1504.01942.

100.  Karpathy, A.; Toderici, G.; Shetty, S.; Leung, T.; Fei-Fei, L. Large-Scale Video Classification with Convolutional Neural Networks. In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, Columbus, OH, USA, 23–28 June 2014.

101.  Kuehne, H.; Jhuang, H.; Stiefelhagen, R.; Serre, T. HMDB: A Large Video Database for Human Motion Recognition. In Proceedings of the 2011 International Conference on Computer Vision, Barcelona, Spain, 6–13 November 2011.

102.  Awad, G.; Butt, A.; Fiscus, J.; Joy, D.; Huet, B. Trecvid 2017: Evaluating ad-hoc and instance video search, events detection, video captioning and hyperlinking. In Proceedings of the TRECVID 2017, Gaithersburg, MD, USA, 13–15 November 2017; Available online: https://hal.archives-ouvertes.fr/hal-01854790 (accessed on 20 August 2020).

103.  Dollar, P.; Wojek, C.; Schiele, B.; Perona, P. Pedestrian detection: A benchmark. In Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition, Miami, FL, USA, 20–25 June 2009; pp. 304–311.

104.  Everingham, M.; Eslami, S.M.A.; Gool, L.V.; Williams, C.K.I.; Winn, J.; Zisserman, A. The Pascal Visual Object Classes Challenge: A Retrospective. *Int. J. Comput. Vis.* **2015**, *111*, 98–136. [CrossRef]

105.  Everingham, M.; Gool, L.V.; Williams, C.K.I.; Winn, J.; Zisserman, A. The Pascal Visual Object Classes (VOC) Challenge. *Int. J. Comput. Vis.* **2010**, *88*, 303–338. [CrossRef]

106.  Han, G.; Zhang, X.; Li, C. Semi-Supervised DFF: Decoupling Detection and Feature Flow for Video Object Detectors. In Proceedings of the 26th ACM international conference on Multimedia, Seoul, Korea, 22–26 October 2018; pp. 1811–1819.

107.  Yang, Y.; Shu, G.; Shah, M. Semi-supervised Learning of Feature Hierarchies for Object Detection in a Video. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, 23–28 June 2013; pp. 1650–1657.

108.  Kumar Singh, K.; Xiao, F.; Jae Lee, Y. Track and Transfer: Watching Videos to Simulate Strong Human Supervision for Weakly-Supervised Object Detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 3548–3556.

109.  Sharma, P.; Huang, C.; Nevatia, R. Unsupervised Incremental Learning for Improved Object Detection in a Video. In Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, 16–21 June 2012; pp. 3298–3305.

110.  Mao, H.; Yang, X.; Dally, W.J. A Delay Metric for Video Object Detection: What Average Precision Fails to Tell. In Proceedings of the IEEE International Conference on Computer Vision, Seoul, Korea, 27 October–2 November 2019.

111.  Dosovitskiy, A.; Fischer, P.; Ilg, E.; Haeusser, P.; Hazirbas, C.; Golkov, V.; van der Smagt, P.; Cremers, D.; Brox, T.; IEEE. FlowNet: Learning Optical Flow with Convolutional Networks. In Proceedings of the 2015 Ieee International Conference on Computer Vision, Santiago, Chile, 7–13 December 2015; pp. 2758–2766.

112. He, K.; Zhang, X.; Ren, S.; Sun, J.; IEEE. Deep Residual Learning for Image Recognition. In Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 770–778.

113. Luo, C.; Zhan, J.; Wang, L.; Yang, Q. Cosine Normalization: Using Cosine Similarity Instead of Dot Product in Neural Networks. In *International Conference on Artificial Neural Networks*; Springer: Cham, Switzerland, 2017.

114. Deng, J.; Zhou, Y.; Yu, B.; Chen, Z.; Zafeiriou, S.; Tao, D. Speed/Accuracy Tradeoffs for Object Detection From Video. Available online: http://image-net.org/challenges/talks_2017/Imagenet%202017%20VID.pdf (accessed on 20 August 2020).

115. Howard, A.G.; Zhu, M.; Chen, B.; Kalenichenko, D.; Wang, W.; Weyand, T.; Andreetto, M.; Adam, H. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. *arXiv* **2017**, arXiv:1704.04861.

116. Hochreiter, S.; Schmidhuber, J. Long short-term memory. *Neural Comput.* **1997**, *9*, 1735–1780. [CrossRef] [PubMed]

117. Lipton, Z.C.; Berkowitz, J.; Elkan, C. A Critical Review of Recurrent Neural Networks for Sequence Learning. *arXiv* **2015**, arXiv:1506.00019.

118. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, L.; Polosukhin, I. Attention Is All You Need. In *Advances in Neural Information Processing Systems 30*; Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R., Eds.; NIPS Proceedings: Denver, CO, USA, 2017; Available online: https://papers.nips.cc/book/advances-in-neural-information-processing-systems-30-2017 (accessed on 20 August 2020).

119. Bahdanau, D.; Cho, K.; Bengio, Y. Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv* **2014**, arXiv:1409.0473.

120. Xie, S.; Girshick, R.; Dollár, P.; Tu, Z.; He, K. Aggregated Residual Transformations for Deep Neural Networks. In Proceedings of the 30th IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 5987–5995.

121. Chen, X.; Yu, J.; Wu, Z. Temporally Identity-Aware SSD With Attentional LSTM. *IEEE Trans. Cybern.* **2020**, *50*, 2674–2686. [CrossRef]

122. Chen, X.; Wu, Z.; Yu, J. TSSD: Temporal Single-Shot Detector Based on Attention and LSTM. In Proceedings of the 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems, Madrid, Spain, 1–5 October 2018; Maciejewski, A.A., Ed.; IEEE: Piscataway, NJ, USA, 2018; pp. 5758–5763.

123. Zhu, H.; Wei, H.; Li, B.; Yuan, X.; Kehtarnavaz, N. Real-Time Moving Object Detection in High—Resolution Video Sensing. *Sensors* **2020**, *20*, 3591. [CrossRef]

124. Zhang, X.; Zhou, X.; Lin, M.; Sun, J. ShuffleNet: An Extremely Efficient Convolutional Neural Network for Mobile Devices. In Proceedings of the 2018 IEEE/Cvf Conference on Computer Vision and Pattern Recognition (CVPR), Salt Lake City, UT, USA, 18–23 June 2018.

125. Liu, Z.; Li, J.; Shen, Z.; Huang, G.; Yan, S.; Zhang, C. Learning Efficient Convolutional Networks through Network Slimming. In Proceedings of the IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017.

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.