

1 *Review*

2 **A review of video object detection: datasets, metrics** 3 **and methods**

4 **Haidi Zhu** ^{1,2}, **Haoran Wei** ³, **Baoqing Li** ^{1,*}, **Xiaobing Yuan** ¹ and **Nasser Kehtarnavaz** ³

5 ¹ Science and Technology on Micro-system Laboratory, Shanghai Institute of Microsystem and Information
6 Technology, Chinese Academy of Sciences, Shanghai 201800, China; hdzhu@mail.sim.ac.cn (H.Z.);
7 sinowsn@mail.sim.ac.cn (X.Y.)

8 ² University of Chinese Academy of Sciences, Beijing 100049, China

9 ³ Department of Electrical and Computer Engineering, University of Texas at Dallas, Richardson, TX 75080,
10 USA; Haoran.Wei@utdallas.edu (H.W.); kehtar@utdallas.edu (N.K.)

11 * Correspondence: sinoiot@mail.sim.ac.cn

12 Received: date; Accepted: date; Published: date

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

26 **Keywords:** video object detection; review of video object detection; deep learning-based video
27 object detection
28

29 **1. Introduction**

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

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

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

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

50 Initially, video object detection approaches relied on handcrafted features, e.g. [29-42]. With the
51 rapid development of deep learning and convolutional neural networks, deep learning models have
52 shown to be more effective than conventional approaches for various tasks in computer vision [43-
53 50], speech processing [51-55], and multi-modality signal processing [56-61]. A number of deep
54 learning-based video object detection approaches were developed after the ILSVRC2015 challenge.
55 These approaches can be divided into flow based [22, 27, 28, 62-64], LSTM based [65-68], attention
56 based [25, 69-72], tracking based [26, 73-77] and other methods [36, 78-85]. A review of these
57 approaches is provided in this paper.

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

62 2. Datasets and Evaluation Metrics

63 2.1. Datasets

64 The most commonly used dataset is the ImageNet VID dataset [5], which is a prevalent
65 benchmark for video object detection. The dataset is split into a training set and a validation set,
66 containing 3862 video snippets and 555 video snippets, respectively. The video streams are annotated
67 on each frame at the frame rate of 25 or 30 fps. In addition, this dataset contains 30 object categories,
68 which are a subset of the categories in the ImageNet DET dataset [86].

69 In the ImageNet VID dataset, the number of objects in each frame is small compared with the
70 datasets used for static image object detection such as COCO [87]. Though the ImageNet VID dataset
71 is widely used, it has limitations in fully reflecting the effect of various video object detection methods.
72 In [88], a large-scale dataset named YouTube-BoundingBoxes (YT-BB) was provided which is human-
73 annotated at one frame per second on video snippets from YouTube with high accuracy classification
74 labels and tight bounding boxes. YT-BB contains approximately 380,000 video segments with 5.6
75 million bounding boxes of 23 object categories which is a subset of the COCO label set. However, the
76 dataset contains only 23 object categories and the image quality is relatively low due to its collection
77 by hand-held mobile phones.

78 In 2018, a dataset named EPIC KITCHENS was provided in [89], which consists of 32 different
79 kitchens in 4 cities with 11,500,000 frames containing 454,158 bounding boxes spanning 290 classes.
80 However, its kitchen scenario poses limitation for performing generic video object detection. Also,
81 there exist the following other datasets that reflect specific applications: the DAVIS dataset [90] for
82 object segmentation, CDnet2014 [91] for moving object detection, VOT [92] and MOT [93] for object
83 tracking. In addition, some works based on semi-supervised or unsupervised methods have been
84 considered in [94-97].

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

89 2.2. Evaluation Metrics

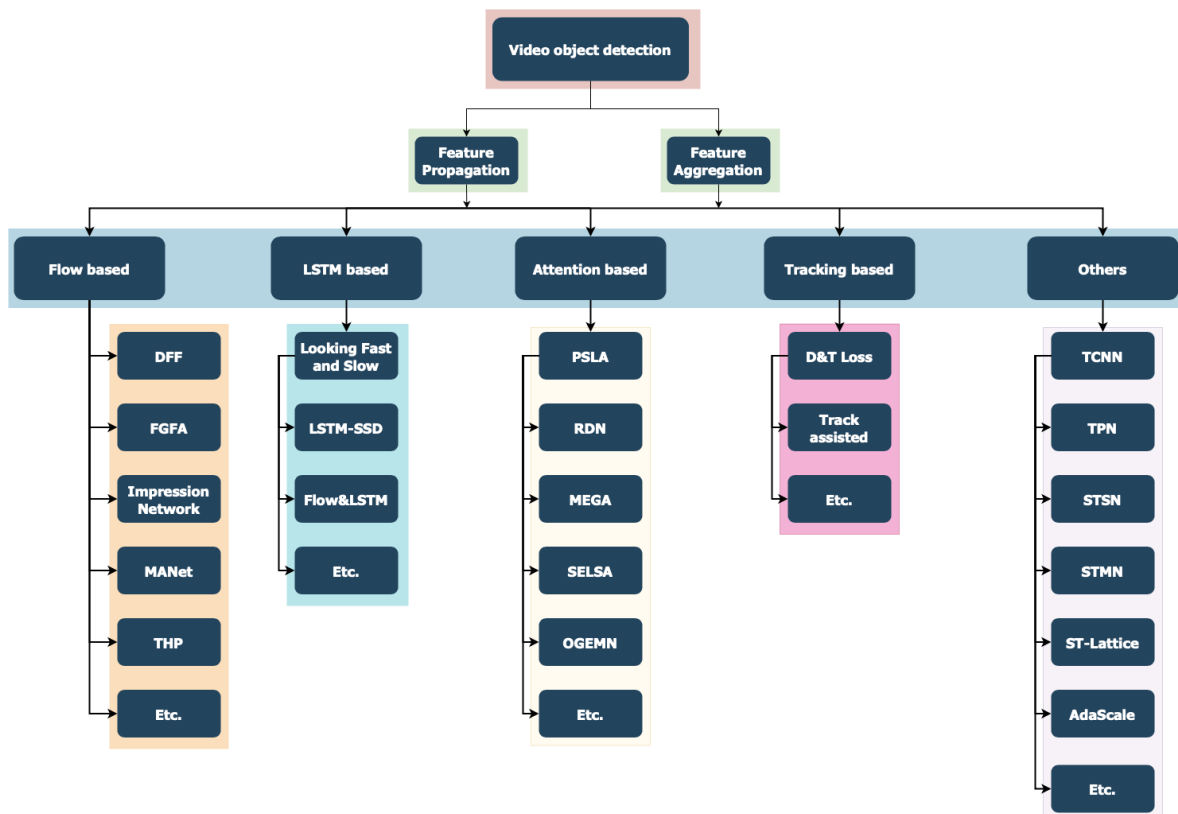
90 The metric mean Average Precision (mAP) is extensively used in conventional object detection,
91 which provides a performance evaluation in terms of regression and classification accuracies [9-15,
92 17]. For video object detection, mAP is also directly used as an evaluation metric in [22, 25, 28, 67, 69].
93 Based on the object speed, it is labeled as mAP (slow), mAP (medium), and mAP (fast) [22]. This is

94 done using the average score of IoU (Intersection over Union) of a current frame and 10 frames ahead
 95 and past as follows: slow (score > 0.9), medium (score \in [0.7, 0.9]), and fast (score < 0.7).

96 In [98], it was pointed out that performance cannot be sufficiently evaluated using only Average
 97 Precision (AP) since the temporal nature of video snippets do not get captured by it. In the same
 98 paper, a new metric named Average Delay (AD) was introduced based on the number of frames
 99 taken to detect an object starting from the frame it first appears. A subset of the ImageNet VID dataset,
 100 named ImageNet VIDT, was considered to verify the effectiveness of AD. It was reported that most
 101 methods having higher ADs still had good APs, indicating that AP was not sufficient to reflect the
 102 temporal characteristics of video object detectors.

103 **3. Video Object Detection Methods**

104 For video object detection, in order to make full use of the video characteristics, different
 105 methods are considered to capture the temporal-spatial relationship. Some papers have considered
 106 the traditional methods [29-42]. These papers heavily rely on the manual design leading to the
 107 shortcomings of low accuracy and lack of robustness to noise sources. More recently, deep learning
 108 solutions have attempted to overcome these shortcomings. As shown in Figure 1, based on the
 109 utilization of the temporal information and the aggregation of features extracted from video snippets,
 110 video object detectors can be divided into flow based [22, 27, 28, 62-64], LSTM based [65-68], attention
 111 based [25, 69-72], tracking based [26, 73-77] and other methods [36, 78-85]. These methods are
 112 described in more detail below.



113

114

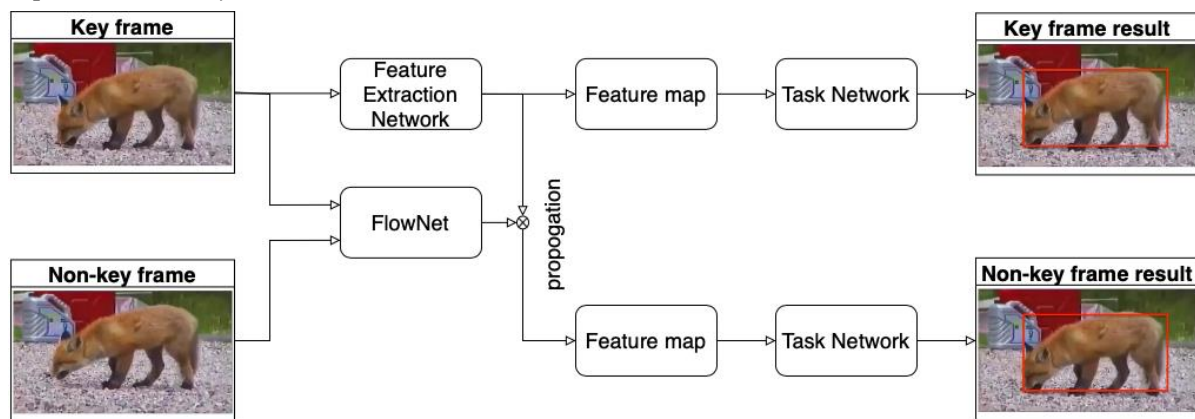
Figure 1. Categories of video object detection methods.

115 *3.1. Flow Based*

116 Flow based methods use optical flow in two ways. In order to save computation, in the first way
 117 as discussed in [28] (DFF), optical flow is used to propagate features from key frames to non-key
 118 frames. In the second way, as discussed in [22] (FGFA), optical flow is used to make use of the
 119 temporal-spatial information between adjacent frames to enhance the features of each frame. In the
 120 second way, higher detection accuracies but lower speeds are reported. As a result, attempts were

121 made to combine both of these ways in [63] (Impression Network) and [64] (THP). To obtain the
 122 difference between adjacent frames and utilize the temporal-spatial information at the pixel level, an
 123 optical flow algorithm was proposed in [29]. In [99], the optical flow estimation was achieved by
 124 using the deep learning model of FlowNet.

125 For video object detection, it is challenging to apply the state-of-the-art object detection
 126 approaches for still images directly to each image frame in video data for the reasons stated earlier.
 127 Therefore, based on FlowNet, the DFF method was proposed in [28] to address these shortcomings:
 128 (i) computation time of feature map extraction for each frame in video, (ii) similarity of features
 129 obtained on two adjacent frames, (iii) propagation of feature maps from one frame to another. In [28],
 130 a convolutional neural sub-network, ResNet-101, was employed to extract the feature map on sparse
 131 key frames. Features on non-key frames were obtained by warping the feature map on key frames
 132 with the flow field generated by FlowNet [99] instead of getting extracted by ResNet-101. The
 133 framework is shown in Figure 2. This method accelerates the object detection on non-key frames. On
 134 the ImageNet VID dataset [5], DFF achieved an accuracy of 73.1% mAP with 20 fps while the baseline
 135 accuracy on a single frame was 73.9% with 4 fps. This method significantly advanced the practical
 136 aspect of video object detection.

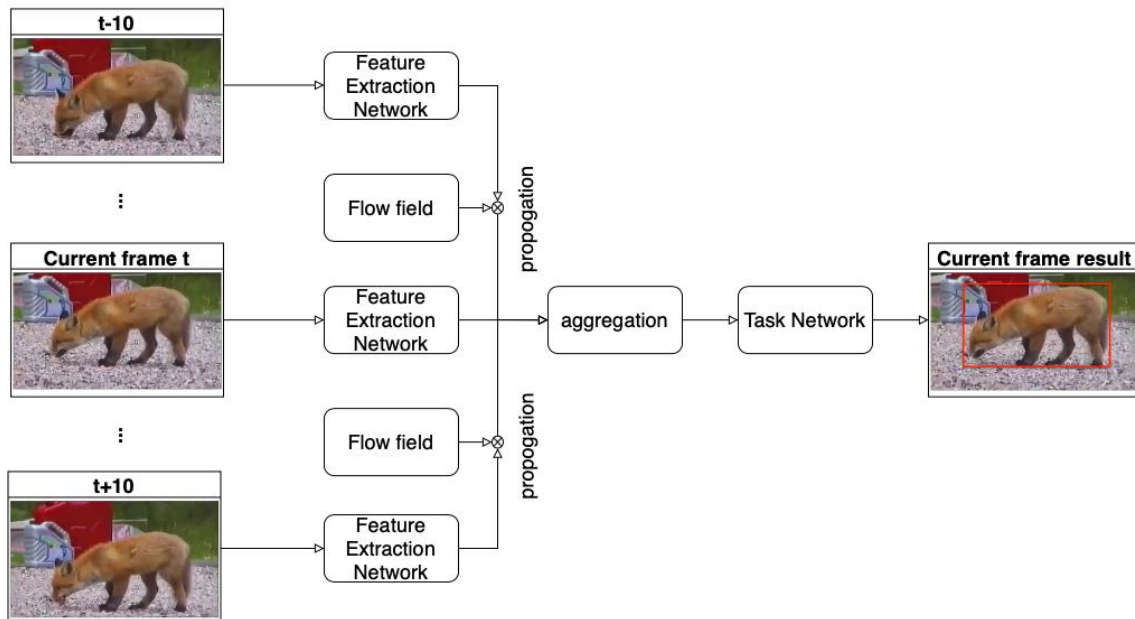


137

138

Figure 2. DFF framework [28].

139 In [22], a flow guided feature aggregation (FGFA) method was proposed to improve the
 140 detection accuracy due to motion blur, rare poses, video defocus, etc. Feature maps were extracted
 141 on each frame in video using ResNet-101 [100]. In order to enhance the feature maps of a current
 142 frame, the feature maps of its nearby frames were warped to the current frame according to the
 143 motion information obtained by the optical flow network. The warped feature maps and extracted
 144 feature maps on the current frame were then inputted into a small sub-network to obtain a new
 145 embedding feature which was used for a similarity measure based on the cosine similarity metric
 146 [101] to compute the weights. Next, the features were aggregated according to the weights. Finally,
 147 the aggregated feature maps were inputted into a shallow detection specific sub-network to obtain
 148 the final detection outcome on the current frame. The framework of FGFA is shown in Figure 3. Based
 149 on the ImageNet VID dataset, FGFA achieved an accuracy of 76.3% mAP with 1.36 fps, which was
 150 higher than DFF.



151

152

Figure 3. FGFA framework [22].

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

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 [63] 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 assistance 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]. Impression Network achieved 75.5% mAP accuracy at 20 fps on the ImageNet VID dataset.

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

Besides Impression Network, in [64] 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 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.

174

175

176

177

178

179

According to the results reported in [64], the mAP accuracy was 78.6% with a runtime of 13.0 and 8.6 fps on the GPUs Titan X and K40, 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 [102] 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 [64].

180

3.2. LSTM Based

181

182

183

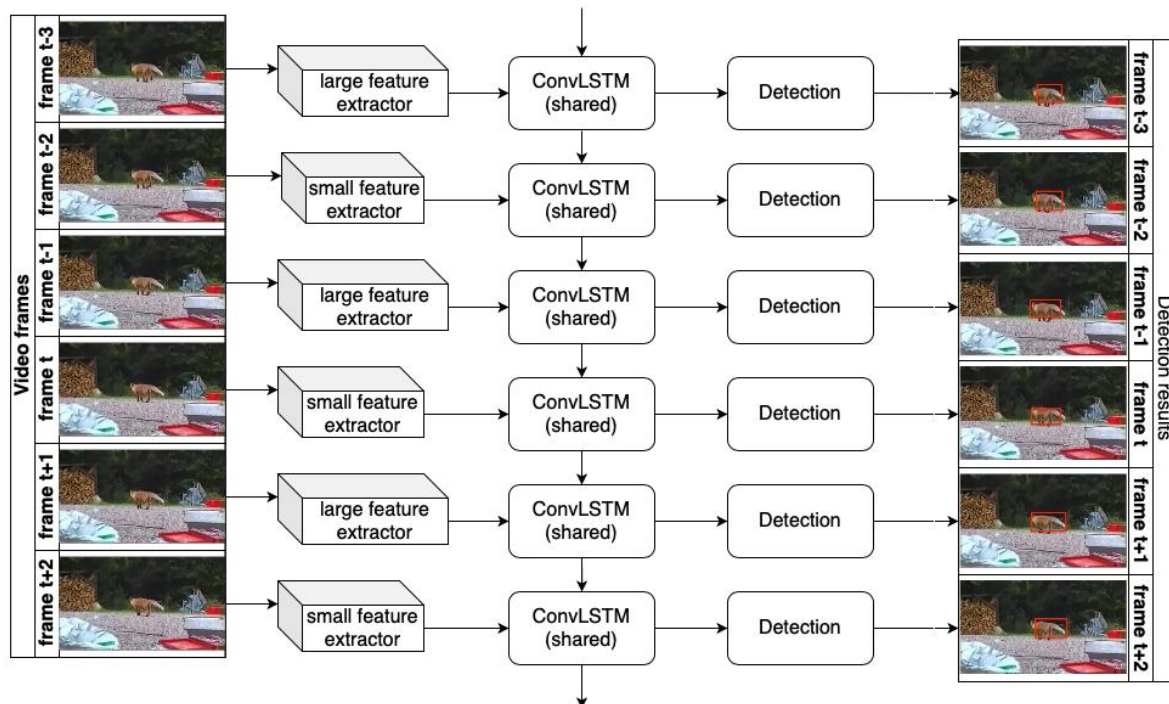
In order to make full use of the temporal-spatial information, convolutional long short term memory (LSTM [103]) was employed to process sequential data in [104] and select important information for a long duration. The methods reported in [65] and [66] are offline LSTM based

184 solutions which utilize all the frames in video. While the method in [67] is an online solution, it only
 185 uses the current and previous frames.

186 In [66], a light model was proposed, which was designed to work on mobile phones and
 187 embedded devices. This method integrated SSD [9] (an efficient object detector network) with the
 188 convolutional LSTM by applying image-based object detector to video object detection via a
 189 convolutional LSTM. The convolutional LSTM was a modified version of the traditional LSTM
 190 encoding the temporal and spatial information.

191 Considering a video snippet as video frames $V = \{I_0, I_1, I_2, \dots, I_t\}$, the model is viewed as a function
 192 $F(I_t, \mathcal{S}_{t-1}) = (D_t, \mathcal{S}_t)$, where D_t denotes the detection outcome of the video object detector and \mathcal{S}_t
 193 represents a vector of feature maps up to the video frame t . Each feature map of \mathcal{S}_{t-1} is the state
 194 input to the LSTM and \mathcal{S}_t is the state output. The state unit \mathcal{S}_t of LSTM contains the temporal
 195 information. LSTM can combine the state unit with input features, adaptively adding the temporal
 196 information to the input features, and updating the state unit at the same time. In [66], it was stated
 197 that such a convolutional LSTM layer could be added to any layer of the original object detector to
 198 refine the input features of the next layer. An LSTM layer could be placed immediately after any
 199 feature map. Placing the LSTM earlier would lead to larger input volumes and much higher
 200 computational cost. In [66], the convolutional LSTM was placed only after the Conv13 layer which
 201 was proved to be most effective through experimental analysis. This method was evaluated on the
 202 ImageNet VID 2015 dataset [5] and achieved a good performance in terms of the model size and
 203 computational efficiency (15 fps on a mobile CPU) with accuracy comparable to those more
 204 computationally demanding single frame models.

205 In 2019, the method in [66] was improved in [65] in terms of inference speed. Specifically, as
 206 shown in Figure 4, due to the high temporal redundancy in video, the model proposed in [65]
 207 contained two feature extractors: a small feature extractor and a large feature extractor. The large
 208 feature extractor with low speed was responsible for extracting the features with high accuracy while
 209 the small feature extractor with fast speed was responsible for extracting the features with poor
 210 accuracy. The two feature extractors were used alternately. The feature maps were aggregated using
 211 a memory mechanism with the modified convolutional LSTM layer. Then, a SSD-style [9] detector
 212 was applied to the refined features to obtain the final regression and classification outcome.



213

214

Figure 4. Small and large feature extractors in [65].

215 For the methods mentioned above, image object detectors together with a temporal context
216 information enhancement were employed to detect objects in video. However, for online video object
217 detection, succeeding frames cannot be utilized. In other words, non-causal video object detectors are
218 not feasible for online applications. Noting that most video object detectors are non-causal, a causal
219 recurrent method was proposed in [67] for online detection without using succeeding frames. In this
220 case, the challenges in terms of occlusion and motion blur remain which require the use of temporal
221 information. For online video object detection, only the current frame and the previous frame are
222 used. Based on the optical flow method [99], the short-term temporal information was utilized by
223 warping the feature maps from the previous frame. However, sometimes image distortion or
224 occlusion would last for several video frames. By using only the short-term temporal information, it
225 was difficult to deal with these situations. The long-term temporal context information was also
226 exploited via the convolutional LSTM, in which the feature maps of the distant preceding frame
227 obtained from the memory function were propagated to acquire more information. Then, the feature
228 maps extracted on the current frame as well as the warped feature maps and the output of the LSTM
229 were concatenated to obtain the aggregated feature maps. Finally, the aggregated feature maps were
230 inputted into a detection sub-network to obtain the detection outcome on the current frame. By
231 utilizing both the short and long-term information, this method achieved an accuracy of 75.5% mAP
232 at a high speed on the ImageNet VID dataset, indicating a competitive performance for online
233 detection.

234 3.3. Attention Related

235 For video object detection, it is known that exploiting the temporal context relationship is quite
236 important. This relationship needs to be established based on a long-duration video, which requires
237 a large amount of memory and computational resources. In order to decrease the computational
238 resources, an attention mechanism was introduced for feature maps alignment. This mechanism was
239 first proposed for machine translation in [105, 106] and then applied to video object detection in [25,
240 69-72].

241 Some methods only take the global or local temporal information into consideration. Specifically,
242 the method RDN in [70] only makes use of the local temporal information. The methods SELSA in
243 [72], OGEMN in [69] only utilize the global temporal information. While the other methods of PSLA
244 in [71], MEGA in [25] use both the global and local temporal information.

245 Relation Distillation Networks (RDN) presented in [70] propagate and aggregate the feature
246 maps based on object relationships in video. In RDN, ResNet-101 [100] and ResNeXt-101-64×4d [107]
247 are utilized as the backbone to extract feature maps and object proposals are generated with the help
248 of a Region Proposal Network (RPN) [15]. The feature maps of each proposal on the reference frame
249 are augmented on the basis of supportive proposals. A prominent innovation in this work is to distill
250 the relation with multi-stage reasoning consisting of a basic and an advanced stage. In the basic stage,
251 the supportive proposals consisting of Top K proposals of a current frame and its adjacent frames are
252 used to measure the relation feature of each reference proposal obtained on the current frame to
253 generate refined reference proposals. In the advanced stage, supportive proposals with high objective
254 scores are selected to generate advanced supportive proposals. Features of selected supportive
255 proposals are aggregated with the relation against all supportive proposals. Then, such aggregated
256 features are employed to strengthen reference proposals obtained from the basic stage. Finally, the
257 aggregated features of reference proposals obtained from the advanced stage are used to generate
258 the final classification and bounding box regression. In addition, the detection box linking is used in
259 a post-processing stage to refine the detection outcome. Evaluated on the ImageNet VID dataset,
260 RDN achieved a detection accuracy of 81.8% and 83.2% mAP, respectively, with ResNet-101 and
261 ResNeXt-101 for feature extraction. With linking and rescore operations, it achieved an accuracy of
262 83.8% and 84.7% mAP, respectively.

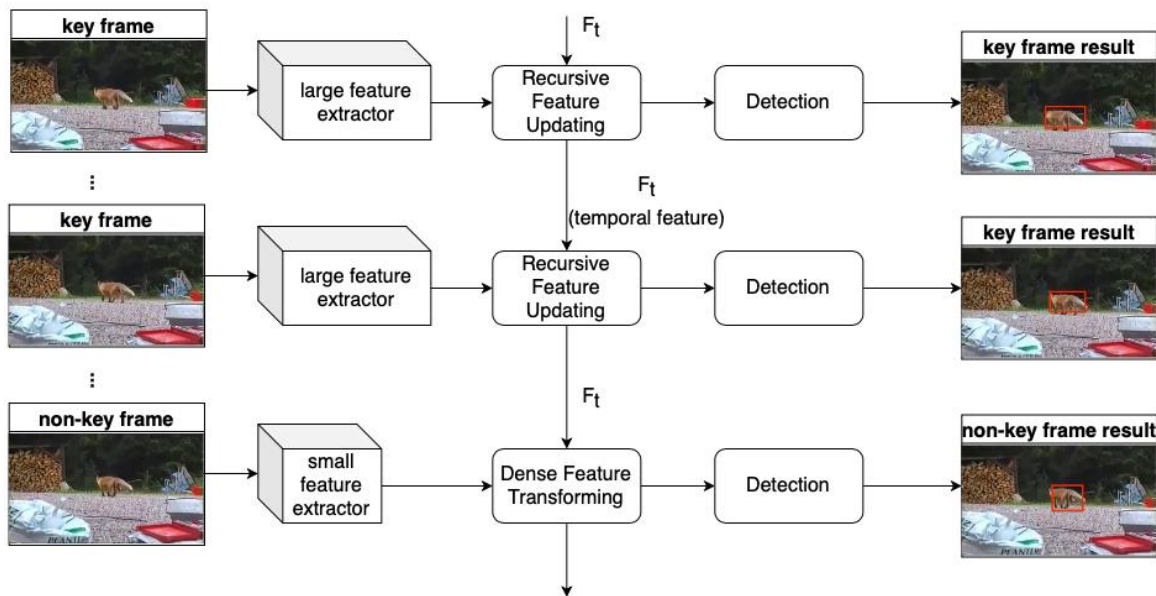
263 A module (SELSA) was introduced in [72] to exploit the relation between the proposals in the
264 entire sequence level, then related feature maps were fused for classification and regression. More
265 specifically, the features of the proposals were extracted on different frames and then a clustering

266 module and a transformation module were applied. The similarities of the proposals were computed
267 across frames and the features were aggregated according to the similarities. Consequently, more
268 robust features were generated for the final detection.

269 In [69], OGEMN was presented which used an object guided external memory to store the pixel
270 and instance level features for further global aggregation. In order to improve the storage-efficiency
271 aspect, only the features within the bounding boxes were stored for further feature aggregation.

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

290 The method Progressive Sparse Local Attention (PSLA) was proposed in [71] to make use of the
291 long term temporal information for enhancement on each feature cell in an attention manner. PSLA
292 establishes correspondence by propagating features in a local region with gradually sparser stride
293 according to the spatial information across frames. Recursive Feature Updating (RFU) and Dense
294 Feature Transforming (DenseFT) were also proposed based on PSLA to model the temporal
295 relationship and enhance the features in a framework shown in Figure 5. More specifically, features
296 were propagated in an attention manner. First, the correspondence between each feature cell in an
297 embedding feature map of a current frame and its surrounding cells was established with a
298 progressive sparser stride from the center to the outside of another embedding feature map of a
299 support frame. Second, correspondence weights were used to compute the aligned feature maps. The
300 feature maps were aggregated with the aligned features. In addition, similar to other video object
301 detectors, the features of key frames were propagated to non-key frames. A light weight network was
302 then applied to extract low-level features on non-key frames and fuse with the features propagated
303 from key frames (DenseFT). Feature propagation was also employed between key frames, and key
304 frame features were updated recursively by an update network (RFU). Hence, features were enriched
305 by the temporal information with DenseFT and RFU, which were further used for detection. Based
306 on the experimentations done in [71], an accuracy of 81.4% mAP was achieved on the ImageNet VID
307 dataset.



308
309

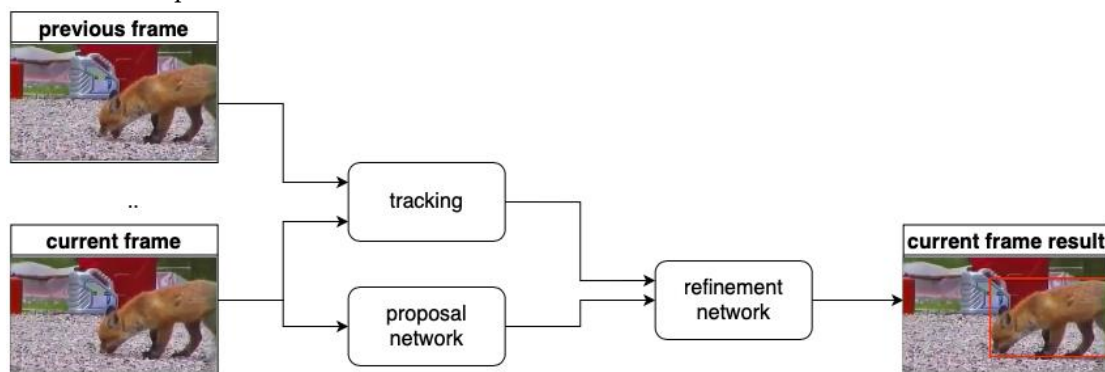
Figure 5. PSLA framework [71].

310 3.4. Tracking Based

311 Inspired by the fact that tracking is an efficient way to utilize the temporal information, several
312 methods [73, 74, 76] have been developed to detect objects on fixed interval frames and track them in
313 frames in between. The improved methods in [26] and [75] detect interval frames adaptively and
314 track the other frames.

315 A framework named CDT was presented in [74] combining detection and tracking for video
316 object detection. This framework consisted of an object detector, a forward tracker and a backward
317 tracker. Initially, objects were detected by the image object detector. Then, each detected object was
318 tracked by the forward tracker, and undetected objects were stored by the backward tracker. In the
319 entire process, the object detector and the tracker cooperated with each other to deal with the
320 appearance and disappearance of objects.

321 Another framework named CaTDet having high computational efficiency was presented in [73].
322 This framework is shown in Figure 6, which includes a tracker and a detector. CaTDet uses a tracker
323 to predict the position of objects with high confidence in a next frame. The processing steps of CaTDet
324 are: (i) Every frame is inputted to a proposal network to output potential proposals in the frame. (ii)
325 Object position in a next frame is predicted with a high confidence using the tracker. (iii) In order to
326 obtain the calibrated object information, the outputs of the tracker and the proposal network are
327 combined and inputted to a refinement network.



328
329

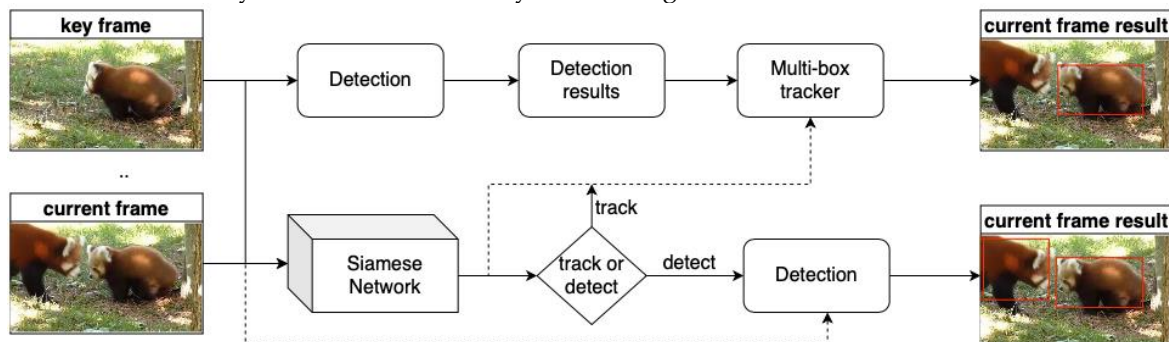
Figure 6. CaTDet framework [73].

330 More specifically, based on the observation that objects detected in one video frame would most
331 likely appear in a next frame, a tracker was used to predict the positions on the next frame with the

332 historical information. In case new objects appeared in a current frame, a computationally efficient
 333 proposal network similar to RPN was utilized to detect proposals. In addition, to address situations
 334 such as motion blur and occlusion, the temporal information was used by a tracker to predict future
 335 positions. The results obtained by combining the tracker and the proposal network was then refined
 336 by a refinement network. Only the regions of interest were refined by the refinement network to save
 337 computation time while maintaining accuracy.

338 Similar to CDT and CaTDet, recent approaches for detection and tracking of objects in video
 339 involve rather complex multistage components. In [76], a framework using a ConvNet architecture
 340 was deployed in a simple but effective way by performing tracking and detection simultaneously.
 341 More specifically, first R-FCN [19] was employed to extract the feature maps shared between
 342 detection and tracking. Then, proposals in each frame were obtained by using RPN based on anchors
 343 [15]. RoI pooling [15] was utilized for the final detection. In particular, a regressor was introduced to
 344 extend the architecture. Position-sensitive regression maps from both frames were used together with
 345 correlation maps as the input to an RoI tracking module, in which the box relationship between the
 346 two frames was outputted. For video object detection, the framework in [76] was evaluated on the
 347 ImageNet VID dataset achieving an accuracy of 82.0% mAP.

348 Similarly, inspired by the observation that object tracking is more efficient than object detection,
 349 a framework (D or T) was covered in [75], see Figure 7, which includes a scheduler network to
 350 determine the operation (detecting or tracking) on a certain frame. Compared with the baseline frame
 351 skipping (detecting on fixed interval frames and tracking on intermediate frames), the scheduler
 352 network with light weights and a simple structure was found to be more effective on the ImageNet
 353 VID dataset. Also, the adaptive mechanism in [26] (TRACKING ASSISTED) was used to select key
 354 frames. Detection on key frames involved the utilization of an accurate detection network and
 355 detection on non-key frames was assisted by the tracking module.



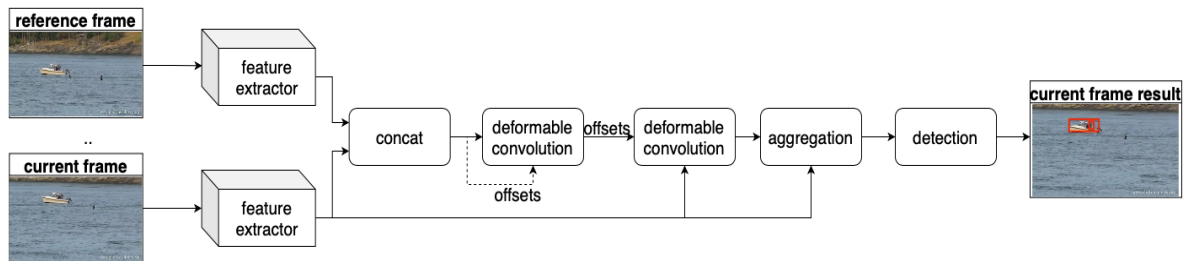
356
357
358 **Figure 7.** D or T framework [75].

358 3.5. Other Methods

359 Apart from the frameworks described above, some methods are presented that are based on a
 360 combination of multiple methods described above [24, 108, 109]. The method in [24] is based on the
 361 optical flow and tracking methods. The methods in [108] (Attentional LSTM) and [109] (TSSD) are
 362 based on the attention and LSTM methods.

363 In addition, these other methods appear in the literature [36, 78-85]. The methods in [78] and [82]
 364 discuss ways to align and enhance feature maps. While the method in [85] studied the effect of the
 365 input image size by selecting a size to achieve a better speed-accuracy trade-off. The method in [78]
 366 named STSN is shown in Figure 8. This method aligns feature maps between adjacent frames. Similar
 367 to the FGFA method in [22], it relies on the idea that detection on a single frame would have
 368 difficulties dealing with noise sources such as motion blur and video defocus. Multiple frames are
 369 thus utilized for feature enhancement to achieve better performance. Unlike FGFA which uses the
 370 optical flow method to align feature maps, deformable convolution is employed for feature alignment
 371 in [78]. First, a sharing feature extraction network is applied to extract feature maps on a current
 372 frame and adjacent frames. Then, the two feature maps are concatenated per channel and a
 373 deformable convolution is performed. The result of the deformable convolution is used as the offset

374 for the second deformable convolution operation to align the feature maps. Furthermore, augmented
 375 feature maps are obtained by aggregating the features in the same way as FGFA. Compared with
 376 FGFA, STSN uses deformable convolution to align the features of two adjacent frames implicitly.
 377 Although it is not as intuitive as the optical flow method, it is also found to be effective. According
 378 to the experimental results reported, STSN still achieved a higher mAP than FGFA (78.9% vs 78.8%)
 379 without relying on the optical flow information. In addition, without the assistant of the temporal
 380 post-processing, STSN obtained a better performance than the D&T baseline [76], 78.9% vs. 75.8%.



381

382

Figure 8. STSN framework [78].

383 Different from [78] by using the deformable convolution to propagate the temporal information,
 384 the Spatial-Temporal Memory Network (STMN) was considered in [82], which involved a RNN
 385 architecture with Spatial-Temporal Memory module (STMM) to incorporate the long-term temporal
 386 information. The Spatial-Temporal Memory Network (STMN) operates in an end-to-end manner to
 387 model the long-term information and align the motion dynamics for video object detection. STMM
 388 is the core module in STMN, a convolutional recurrent computation unit which fully utilizes the
 389 pretrained weights learned from static image datasets such as ImageNet [86]. This design is essential
 390 to address the practical difficulties of learning from video datasets, which largely lack the diversity
 391 of objects within the same category. STMM receives the feature maps of a current frame at time step
 392 t and the spatial-temporal memory M_{t-1}^{\rightarrow} with the information of all the previous frames. Then, the
 393 spatial-temporal memory M_t^{\rightarrow} of the current time step is updated. In order to capture the information
 394 of both later frames and previous frames at the same time, two STMMs are used for bidirectional
 395 feature aggregation to produce the memory M which is employed for both classification and
 396 bounding box regression. Therefore, the feature maps are propagated and aggregated by combining
 397 the information across multiple video frames. Evaluated on the ImageNet VID dataset, STMN has
 398 achieved the current start-of-the-art accuracy.

399

All the algorithms described above start from how to propagate and aggregate feature maps. In
 400 [85], video object detection was examined from another point of view. Similar to [110], the effect of
 401 input image size on the performance of video object detection was studied in [85]. Furthermore, it
 402 was found that down sampling images can obtain better accuracy sometime. From this point of view,
 403 a framework named AdaScale was proposed to adaptively select the input image size. AdaScale
 404 predicts the best scale or size of a next frame according to the information of a current frame. One
 405 of the reasons for the improvement is that the number of false positives is reduced. And the other reason
 406 is that the number of true positives is increased by resizing the too large objects to a suitable size for
 407 the detector.

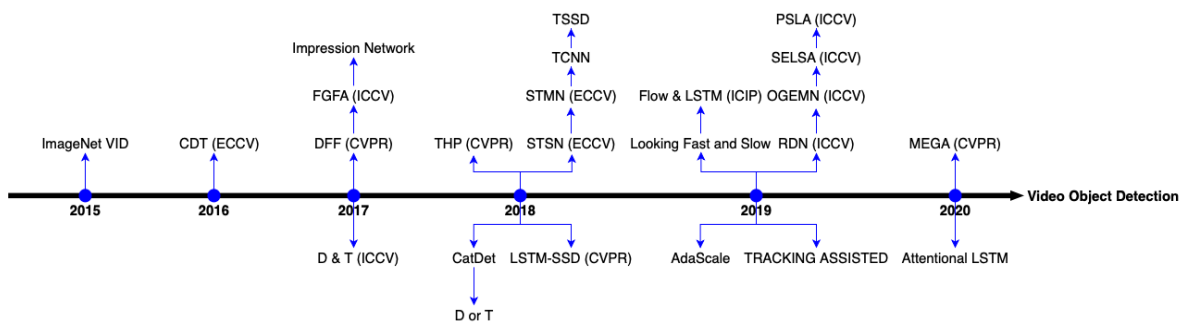
408

In [85], the optimal scale (pixels of the shortest side) of a given image is defined with a predefined
 409 finite set of scales S ($S = \{600, 480, 360, 240\}$ in [85]). Furthermore, a loss function consisting of the
 410 classification and regression loss is employed as the evaluation metric to compare the results across
 411 different scales. The regression loss for background is expected to be zero. Hence, if the loss function
 412 is utilized directly to evaluate the results across different scales, the image scale which contains fewer
 413 foreground bounding boxes is supported. In order to deal with this issue, a new metric (the loss
 414 function which focuses on the same number of foreground bounding boxes chosen on different scales)
 415 is employed to compare across different scales. More specifically, the number of bounding boxes
 416 involved to compute the loss is determined by the minimum number (m) on all the scales. For each
 417 scale, the loss of predicted foreground bounding boxes on the image is sorted in ascending order and

418 the first m bounding boxes are chosen. The scale m with the minimum loss is defined as the best scale.
 419 Inspired by R-FCN[19] working on deep features for bounding boxes regression, the channels of the
 420 deep features are expected to contain the size information. Therefore, a scale regressor using deep
 421 features is built to predict the optimal scale. Evaluated on the ImageNet VID and mini YouTube-BB
 422 datasets, Adascale achieved 1.3% and 2.7% mAP improvements with 1.6 and 1.8 times speedup
 423 compared with a single-scale training and testing, respectively. Furthermore, combined with DFF
 424 [28], the speed was increased by 25% while maintaining mAP on the ImageNet VID dataset.

425 **4. Comparison of Video Object Detection Methods**

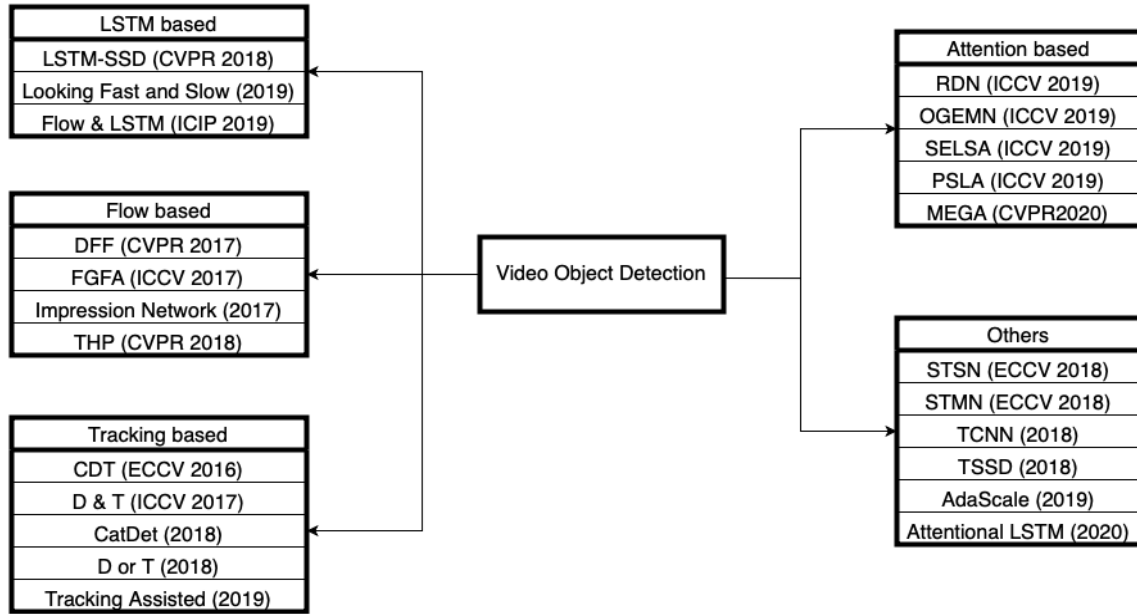
426 The great majority of video object detection approaches use the ImageNet VID dataset [5] for
 427 performance evaluation. In this section, the timeline of video object detection methods in recent years
 428 is shown in Figure 9 together with a group listing of the methods in Figure 10. Then, a comparison is
 429 provided between the methods covered in the previous section. The comparison is presented in Table
 430 1 and Table 2 which correspond to with and without post-processing, respectively. The methods in
 431 Figure 9 belong to different groups but the same time whereas the methods in Figure 10 belong to
 432 different times but the same groups. As can be seen from Figures 9 and 10, the methods based on
 433 optical flow were proposed earlier. During the same period, video object detection methods were
 434 assisted by tracking due to the effectiveness of tracking in utilizing the temporal-spatial information.
 435 The optical flow-based methods needed a large number of parameters and they were only suitable
 436 for small motions. In recent years, the methods based on attention have achieved much success such
 437 as MEGA [25]. Using LSTM for feature propagation and aggregation is becoming a hot research topic
 438 and many new methods are being proposed such as STSN [78] using deformable convolution to align
 439 the feature maps. The latest research is mostly based on attention, LSTM or combination of methods
 440 such as Flow & LSTM [67].



441

442

Figure 9. Timeline of video object detection methods.



443

444

Figure 10. Video object detection methods sorted in different groups.

445

446

447

448

Table 1. Comparison among the video object detection methods without post processing; note that the runtime is based on the 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
Flow Based	Impression Network[63]	Modified ResNet-101	75.5	20 1060
	FGFA [22]	ResNet-101	76.3	1.36 K
	DFF [28]	ResNet-101	73.1	20.25 K
	THP [64]	ResNet-101+DCN	78.6	13.0X/8.6K
LSTM Based	Looking Fast and Slow [65]	Interleaved $\begin{bmatrix} I \\ SEP \end{bmatrix}$	61.4	23.5 Pixel 3 phone
	LSTM-SSD[66]	MobilenetV2-SSDLite	53.5	-
	Flow&LSTM [67]	ResNet-101	75.5	-
Attention Based	OGEMN[69]	ResNet-101	79.3	8.9 (1080Ti)
		ResNet-101+DCN	80.0	-
	PSLA[71]	ResNet-101	77.1	30.8V \ 18.73X
		ResNet-101+DCN	80.0	26.0V \ 13.34X
	SELSA[72]	ResNet-101	80.25	-
		ResNeXt-101	83.11	-
RDN[70]	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[76]	ResNet-101	75.8	7.8X
	Track assisted[26]	ResNet-101	70.0	30XP
Others	TCNN[24]	GoogLeNet	73.8	-
	STSN [78]	ResNet-101+DCN	78.9	-

449

450

451

Table 2. 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	-	
LSTM Based	Looking Fast and Slow[65]	Interleaved		72.3	
		+ Quantization _{SEP} ^[11] + Async _{SEP} ^[11]	59.3	<i>Pixel 3 phone</i>	
		MobilenetV2-SSDLite + LSTM ($\alpha = 1.4$)[66]	MobilenetV2-SSDLite	64.1	4.1 <i>Pixel 3 phone</i>
		MobilenetV2-SSDLite + LSTM($\alpha = 1.0$) [66]	MobilenetV2-SSDLite	59.1	-
		MobilenetV2-SSDLite + LSTM($\alpha = 0.5$) [66]	MobilenetV2-SSDLite	50.3	-
	MobilenetV2-SSDLite + LSTM ($\alpha = 0.35$) [66]	MobilenetV2-SSDLite	45.1	14.6 <i>Pixel 3 phone</i>	
Attention Based	OGEMN [69]	ResNet-101	80.8	-	
		ResNet-101+DCN	81.6	-	
		ResNet-101	78.6	5.7X	
		ResNet-101+DCN	81.4	6.31V\5.13X	
		ResNet-101	80.54	-	
		ResNeXt-101	84.7	-	
Tracking Based	D&T ($\tau = 10$) [76]	ResNet-101	78.6	-	
		ResNet-101	79.8	5X	
		Inception V4	82.0	-	
Others	STSN [78]	ResNet-101+DCN	80.4	-	
		ResNet-101	80.5	-	

452 Table 1 provides the outcomes without post processing. In this table, the methods are divided
453 into different groups according to the way temporal and spatial information are utilized. Flow-
454 guided group propagate and align the feature maps according to the flow field obtained by optical
455 flow. Both accuracy and speed of various frameworks are reported in this table. For example, DFF
456 provides high computational efficiency and achieves a runtime of 20.25 fps using a Titan K40 GPU.
457 FGFA achieves a high accuracy producing 76.3% mAP with 1.36 fps. Obviously, DFF is faster than
458 FGFA. Flow-guided methods are intuitive and well understood to propagate features. Optical flow
459 is deemed suitable for small movement estimation. In addition, since optical flow reflects pixel level
460 displacement, it has difficulties when it is applied to high-level feature maps. One pixel movement
461 on feature maps may correspond to 10 to 20 pixels movement.

462 Inspired by the LSTM based solutions in natural language processing, LSTM methods are used
463 to incorporate the sequence information. In the LSTM group, Flow & LSTM [67] achieved the highest
464 accuracy of 75.5%. Looking Fast and Slow [65] generated high speed but with low accuracy. LSTM
465 captures the long term information with a simple implementation. Since the sigmoid activation of the
466 input and forget gates are rarely completely saturated, a slow state decay and thus loss of long-term
467 dependence is resulted. In other words, it is difficult to retain the complete previous state in the
468 update.

469 Attention based methods also show the ability to perform video object detection effectively. In
470 the attention related group, MEGA [25] with ResNeXt-101 as backbone achieved the highest accuracy
471 of 84.1% mAP. As described, it achieved a very high accuracy with a relatively fast speed. Attention
472 based methods aggregate the features within proposals that are generated. This decreases the
473 computation time. Because of only using the features within the proposals, the performance relies on
474 the effect of RPN to a certain extent. Here, it is rather difficult to utilize more comprehensive
475 information.

476 In the tracking based group, the methods are assisted by tracking. D&T loss [76] achieved 75.8%
477 mAP. Tracking is an efficient method to employ the temporal information with a detector assisted by
478 a tracker. However, it cannot solve the problems created by motion blur and video defocus directly.
479 As the detection performance relies on the tracking performance, the detector part suffers from
480 tracking errors. There are also other standalone methods including TCNN[24], STSN [78] and STMN
481 [82].

482 In order to further improve the performance in terms of detection accuracy, post processing can
483 be added to the above methods. The results with post processing are shown in Table 2. One can easily
484 see that with post processing, the accuracy is noticeably improved. For example, the accuracy of
485 MEGA is improved from 84.1% to 85.4% mAP.

486 5. Future Trends

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

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

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

504 Most of the methods covered in this review paper only utilize the local temporal information or
505 global information separately. There are only a few methods such as MEGA, which have used the
506 local and global temporal information at the same time and achieved a benchmark mAP of 85.4%. As
507 demonstrated by MEGA, it is worth developing future frameworks which utilize both the local and
508 global temporal information. Furthermore, for most of the existing video object detection algorithms,
509 the number of frames used is too small to fully utilize the video information. Hence, as yet another
510 future trend, it is of importance to develop methods that utilize the long-term video information. As
511 can be observed from Tables 1 and 2, the attention-based frameworks achieved a relatively high
512 accuracy. However, such methods pose difficulties for real-time applications demanding very
513 powerful GPUs. Although the Looking Fast and Slow method [65] achieved 72.3 fps on Pixel 3
514 phones, the accuracy is only 59.3% which poses challenging for actual deployment. Indeed, the trade-
515 off between accuracy and speed needs to be further investigated.

516 6. Conclusion

517 In recent years, after the ImageNet Large Scale Visual Recognition Challenge (ILSVRC)
518 announced the video object detection task in 2015, many deep learning-based video object detection
519 solutions have been developed. This paper has provided a review of the video object detection
520 methods that have been developed so far. This review has covered the available datasets, evaluation
521 metrics and an overview of different categories of deep learning-based methods for video object
522 detection. A categorization of the video object detection methods has been made according to the
523 way temporal and spatial information are used. These categories include flow based, LSTM based,
524 attention based, tracking based, as well as other methods. The performance of various detectors with
525 or without post-processing is summarized in two tables in terms of both detection accuracy and

526 computation speed. Several trends of importance in video object detection have also been stated for
527 possible future works.

528 **Author Contributions:** Investigation, H.Z. and H.W.; resources, H.Z., B.L. and X.Y.; writing—original draft
529 preparation, H.Z. and H.W.; rewriting and editing, H.Z., B.L., X.Y. and N.K.; supervision, B.L., X.Y. and N.K. All
530 authors have read and agreed to the published version of the manuscript.

531 **Funding:** This research received no external funding.

532 **Conflicts of Interest:** The authors declare no conflict of interest.

533 References

- 534 1. Bateni, S., et al., *Co-Optimizing Performance and Memory Footprint Via Integrated CPU/GPU Memory*
535 *Management, an Implementation on Autonomous Driving Platform*. 2020.
- 536 2. Lu, J., et al. *A Review on Object Detection Based on Deep Convolutional Neural Networks for Autonomous*
537 *Driving*. in *2019 Chinese Control And Decision Conference (CCDC)*. 2019.
- 538 3. Wei, H., et al., *Deep Learning-Based Person Detection and Classification for Far Field Video Surveillance*.
539 *Proceedings of the 2018 Ieee 13th Dallas Circuits and Systems Conference*. 2018.
- 540 4. Marielet Guillermo, R.R.T., Luigi Carlo De Jesus, Robert Kerwin Dela Cruz Billones, Edwin Sybingco,
541 Elmer P. Dadios, Alexis Fillone, *Detection and Classification of Public Security Threats in the Philippines*
542 *Using Neural Networks*. In *Proceedings of the 13th IEEE Dallas Circuits and Systems Conference, Dallas,*
543 *TX, USA, 2018: p. 1–4*.
- 544 5. Russakovsky, O., et al., *ImageNet Large Scale Visual Recognition Challenge*. *International Journal of*
545 *Computer Vision*, 2015. **115**(3): p. 211-252.
- 546 6. Shen, Z., et al., *DSOD: Learning Deeply Supervised Object Detectors from Scratch*, in *2017 Ieee International*
547 *Conference on Computer Vision*. 2017. p. 1937-1945.
- 548 7. Tian, Z.a.S., Chunhua and Chen, Hao and He, Tong, *FCOS: Fully Convolutional One-Stage Object*
549 *Detection*. 2019 *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- 550 8. Zhao, Q., et al., *M2Det: A Single-Shot Object Detector Based on Multi-Level Feature Pyramid Network*. *Thirty-*
551 *Third Aaai Conference on Artificial Intelligence / Thirty-First Innovative Applications of Artificial*
552 *Intelligence Conference / Ninth Aaai Symposium on Educational Advances in Artificial Intelligence*.
553 2019. 9259-9266.
- 554 9. Liu, W., et al., *SSD: Single Shot MultiBox Detector*, in *Computer Vision - Eccv 2016, Pt I*, B. Leibe, et al.,
555 *Editors*. 2016. p. 21-37.
- 556 10. Redmon, J., A. Farhadi, and Ieee, *YOLO9000: Better, Faster, Stronger*, in *30th Ieee Conference on Computer*
557 *Vision and Pattern Recognition*. 2017. p. 6517-6525.
- 558 11. Redmon, J. and A. Farhadi, *YOLOv3: An Incremental Improvement*. 2018.
- 559 12. Redmon, J., et al., *You Only Look Once: Unified, Real-Time Object Detection*, in *2016 Ieee Conference on*
560 *Computer Vision and Pattern Recognition*. 2016. p. 779-788.
- 561 13. He, K., et al., *Mask R-CNN*, in *2017 Ieee International Conference on Computer Vision*. 2017. p. 2980-2988.
- 562 14. Girshick, R. and Ieee, *Fast R-CNN*, in *2015 Ieee International Conference on Computer Vision*. 2015. p. 1440-
563 1448.
- 564 15. Ren, S., et al., *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks*. *Ieee*
565 *Transactions on Pattern Analysis and Machine Intelligence*, 2017. **39**(6): p. 1137-1149.
- 566 16. Girshick, R., et al., *Rich feature hierarchies for accurate object detection and semantic segmentation*, in *2014*
567 *Ieee Conference on Computer Vision and Pattern Recognition*. 2014. p. 580-587.

- 568 17. Cai, Z., N. Vasconcelos, and Ieee, *Cascade R-CNN: Delving into High Quality Object Detection*, in 2018
569 *Ieee/Cof Conference on Computer Vision and Pattern Recognition*. 2018. p. 6154-6162.
- 570 18. He, K., et al., *Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition*, in *Computer*
571 *Vision - Eccv 2014, Pt Iii*, D. Fleet, et al., Editors. 2014. p. 346-361.
- 572 19. Dai, J., et al., *R-FCN: Object Detection via Region-based Fully Convolutional Networks*, in *Advances in Neural*
573 *Information Processing Systems 29*, D.D. Lee, et al., Editors. 2016.
- 574 20. Shrivastava, A., et al., *Training Region-based Object Detectors with Online Hard Example Mining*, in 2016
575 *Ieee Conference on Computer Vision and Pattern Recognition*. 2016. p. 761-769.
- 576 21. Wei, H. and N. Kehtarnavaz, *Semi-Supervised Faster RCNN-Based Person Detection and Load Classification*
577 *for Far Field Video Surveillance*. 2019.
- 578 22. Zhu, X., et al., *Flow-Guided Feature Aggregation for Video Object Detection*, in 2017 *Ieee International*
579 *Conference on Computer Vision*. 2017. p. 408-417.
- 580 23. Zhang, R., et al., *Video Object Detection by Aggregating Features across Adjacent Frames*, in 2019 *3rd*
581 *International Conference on Machine Vision and Information Technology*. 2019.
- 582 24. Kang, K., et al., *T-CNN: Tubelets With Convolutional Neural Networks for Object Detection From Videos*. Ieee
583 *Transactions on Circuits and Systems for Video Technology*, 2018. **28**(10): p. 2896-2907.
- 584 25. Chen, Y., et al., *Memory Enhanced Global-Local Aggregation for Video Object Detection*. 2020.
- 585 26. Yang, W., et al., *TRACKING ASSISTED FASTER VIDEO OBJECT DETECTION*, in 2019 *Ieee International*
586 *Conference on Multimedia and Expo*. 2019. p. 1750-1755.
- 587 27. Yuan, X.Z.a.J.D.a.X.Z.a.Y.W.a.L., *Towards High Performance Video Object Detection for Mobiles*. arXiv:
588 1804.05830, 2018.
- 589 28. Zhu, X., et al., *Deep Feature Flow for Video Recognition*, in *30th Ieee Conference on Computer Vision and*
590 *Pattern Recognition*. 2017. p. 4141-4150.
- 591 29. Horn, B.K.P. and B.G. Schunck, *DETERMINING OPTICAL-FLOW*. *Artificial Intelligence*, 1981. **17**(1-3):
592 p. 185-203.
- 593 30. Nguyen, H.T., M. Worring, and A. Dev, *Detection of moving objects in video using a robust motion similarity*
594 *measure*. *Ieee Transactions on Image Processing*, 2000. **9**(1): p. 137-141.
- 595 31. Carminati, L., J. Benois-Pineau, and Ieee, *Gaussian mixture classification for moving object detection in video*
596 *surveillance environment*, in 2005 *International Conference on Image Processing*. 2005. p. 3361-3364.
- 597 32. Jayabalan, E. and A. Krishnan, *Object Detection and Tracking in Videos Using Snake and Optical Flow*
598 *Approach*, in *Computer Networks and Information Technologies*, V.V. Das, J. Stephen, and Y. Chaba, Editors.
599 2011. p. 299+.
- 600 33. Jayabalan, E. and A. Krishnan, *Detection and Tracking of Moving Object in Compressed Videos*, in *Computer*
601 *Networks and Information Technologies*, V.V. Das, J. Stephen, and Y. Chaba, Editors. 2011. p. 39+.
- 602 34. Ghosh, A., B.N. Subudhi, and S. Ghosh, *Object Detection From Videos Captured by Moving Camera by Fuzzy*
603 *Edge Incorporated Markov Random Field and Local Histogram Matching*. *Ieee Transactions on Circuits and*
604 *Systems for Video Technology*, 2012. **22**(8): p. 1127-1135.
- 605 35. Guo, C. and H. Gao, *Adaptive graph-cuts algorithm based on higher-order MRF for video moving object*
606 *detection*. *Electronics Letters*, 2012. **48**(7): p. 371-U103.
- 607 36. Guo, C., et al., *An adaptive graph cut algorithm for video moving objects detection*. *Multimedia Tools and*
608 *Applications*, 2014. **72**(3): p. 2633-2652.
- 609 37. Yadav, D.K. and K. Singh, *A combined approach of Kullback-Leibler divergence and background subtraction*
610 *for moving object detection in thermal video*. *Infrared Physics & Technology*, 2016. **76**: p. 21-31.

- 611 38. Oreifej, O., X. Li, and M. Shah, *Simultaneous Video Stabilization and Moving Object Detection in Turbulence*.
612 Ieee Transactions on Pattern Analysis and Machine Intelligence, 2013. **35**(2): p. 450-462.
- 613 39. Nadimi, S. and B. Bhanu, *Physical models for moving shadow and object detection in video*. Ieee Transactions
614 on Pattern Analysis and Machine Intelligence, 2004. **26**(8): p. 1079-1087.
- 615 40. Utsumi, O., et al., *An object detection method for describing soccer games from video*. Ieee International
616 Conference on Multimedia and Expo, Vol I and II, Proceedings. 2002. 45-48.
- 617 41. Hossain, M.J., M.A.A. Dewan, and O. Chae, *Moving object detection for real time video surveillance: An edge
618 based approach*. Ieee Transactions on Communications, 2007. **E90B**(12): p. 3654-3664.
- 619 42. Chiranjeevi, P. and S. Sengupta, *Robust detection of moving objects in video sequences through rough set
620 theory framework*. Image and Vision Computing, 2012. **30**(11): p. 829-842.
- 621 43. Abd Razak, H., et al., *Anomalous Behaviour Detection using Transfer Learning Algorithm of Series and DAG
622 Network*, in *2019 Ieee 9th International Conference on System Engineering and Technology*. 2019. p. 505-509.
- 623 44. Azarang, A., H.E. Manoochehri, and N. Kehtarnavaz, *Convolutional Autoencoder-Based Multispectral
624 Image Fusion*. Ieee Access, 2019. **7**: p. 35673-35683.
- 625 45. Majumder, S., et al. *A deep learning-based smartphone app for real-time detection of five stages of diabetic
626 retinopathy*. in *Real-Time Image Processing and Deep Learning 2020*. 2020.
- 627 46. Wang, Z., et al. *Few-Sample and Adversarial Representation Learning for Continual Stream Mining*. in *WWW
628 '20: The Web Conference 2020*. 2020.
- 629 47. Maor, G., et al., *An FPGA Implementation of Stochastic Computing-based LSTM*, in *2019 Ieee 37th
630 International Conference on Computer Design*. 2019. p. 38-46.
- 631 48. X, C., *Human Pose Estimation and Immediacy Prediction with Deep Learning*. The Chinese University of
632 Hong Kong (Hong Kong), 2017.
- 633 49. Wang, Z., et al., *Metric Learning based Framework for Streaming Classification with Concept Evolution*, in
634 *2019 International Joint Conference on Neural Networks*. 2019.
- 635 50. Li, H., et al., *Multiple Description Coding Based on Convolutional Auto-Encoder*. Ieee Access, 2019. **7**: p.
636 26013-26021.
- 637 51. Siqu Zheng, G.L., Hongbin Suo, Yun Lei, *Autoencoder-Based Semi-Supervised Curriculum Learning for Out-
638 of-Domain Speaker Verification*. INTERSPEECH, 2019: p. 4360-4364.
- 639 52. Wei, H., N. Kehtarnavaz, and Ieee, *Determining Number of Speakers from Single Microphone Speech Signals
640 by Multi-Label Convolutional Neural Network*, in *Iecon 2018 - 44th Annual Conference of the Ieee Industrial
641 Electronics Society*. 2018. p. 2706-2710.
- 642 53. Zhao, Y., et al., *DNN-BASED ENHANCEMENT OF NOISY AND REVERBERANT SPEECH*, in *2016 Ieee
643 International Conference on Acoustics, Speech and Signal Processing Proceedings*. 2016. p. 6525-6529.
- 644 54. Tao, F., et al., *AN ENSEMBLE FRAMEWORK OF VOICE-BASED EMOTION RECOGNITION SYSTEM
645 FOR FILMS AND TV PROGRAMS*. 2018 Ieee International Conference on Acoustics, Speech and Signal
646 Processing. 2018. 6209-6213.
- 647 55. Zhao, Y., et al., *PERCEPTUALLY GUIDED SPEECH ENHANCEMENT USING DEEP NEURAL
648 NETWORKS*. 2018 Ieee International Conference on Acoustics, Speech and Signal Processing. 2018.
649 5074-5078.
- 650 56. Tao, F. and C. Busso. *Aligning Audiovisual Features for Audiovisual Speech Recognition*. in *IEEE
651 International Conference on Multimedia and Expo*. 2018.
- 652 57. Wei, H., P. Chopada, and N. Kehtarnavaz, *C-MHAD: Continuous Multimodal Human Action Dataset of
653 Simultaneous Video and Inertial Sensing*. Sensors, 2020. **20**(10).

- 654 58. Brena, R.F., et al., *Choosing the Best Sensor Fusion Method: A Machine-Learning Approach*. Sensors, 2020.
655 20(8).
- 656 59. Tao, F. and C. Busso, *End-to-End Audiovisual Speech Recognition System with Multitask Learning*. IEEE
657 Transactions on Multimedia, 2020. PP(99): p. 1-1.
- 658 60. Wei, H. and N. Kehtarnavaz, *Simultaneous Utilization of Inertial and Video Sensing for Action Detection and
659 Recognition in Continuous Action Streams*. Ieee Sensors Journal, 2020. 20(11): p. 6055-6063.
- 660 61. Chen, C., R. Jafari, and N. Kehtarnavaz, *A survey of depth and inertial sensor fusion for human action
661 recognition*. Multimedia Tools and Applications, 2017. 76(3): p. 4405-4425.
- 662 62. Wang, S., et al. *Fully Motion-Aware Network for Video Object Detection*. 2018. Cham: Springer International
663 Publishing.
- 664 63. Yan, C.H.a.H.Q.a.S.L.a.J., *Impression Network for Video Object Detection*. arXiv: 1712.05896, 2017.
- 665 64. Zhu, X., et al., *Towards High Performance Video Object Detection*, in *2018 Ieee/Cvf Conference on Computer
666 Vision and Pattern Recognition*. 2018. p. 7210-7218.
- 667 65. Kalenichenko, M.L.a.M.Z.a.M.W.a.Y.L.a.D., *Looking Fast and Slow: Memory-Guided Mobile Video Object
668 Detection*. arXiv: 1903.10172, 2019.
- 669 66. Liu, M., M. Zhu, and Ieee, *Mobile Video Object Detection with Temporally-Aware Feature Maps*, in *2018
670 Ieee/Cvf Conference on Computer Vision and Pattern Recognition*. 2018. p. 5686-5695.
- 671 67. Zhang, C., J. Kim, and Ieee, *MODELING LONG- AND SHORT-TERM TEMPORAL CONTEXT FOR
672 VIDEO OBJECT DETECTION*, in *2019 Ieee International Conference on Image Processing*. 2019. p. 71-75.
- 673 68. Lu, Y., et al., *Online Video Object Detection using Association LSTM*, in *2017 Ieee International Conference on
674 Computer Vision*. 2017. p. 2363-2371.
- 675 69. Deng, H.a.H., Yang and Song, Tao and Zhang, Zongpu and Xue, Zhengui and Ma, Ruhui and
676 Robertson, Neil and Guan, Haibing, *Object Guided External Memory Network for Video Object Detection*.
677 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019: p. 6677-6686.
- 678 70. Deng, J.a.P., Yingwei and Yao, Ting and Zhou, Wengang and Li, Houqiang and Mei, Tao, *Relation
679 Distillation Networks for Video Object Detection*. 2019 IEEE/CVF International Conference on Computer
680 Vision (ICCV), 2019: p. 7022-7031.
- 681 71. Guo, C., et al., *Progressive Sparse Local Attention for Video object detection*. 2019.
- 682 72. Wu, H.a.C., Yuntao and Wang, Naiyan and Zhang, Zhao-Xiang, *Sequence Level Semantics Aggregation
683 for Video Object Detection*. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019.
- 684 73. Mao, H., T. Kong, and W.J. Dally, *CaTDet: Cascaded Tracked Detector for Efficient Object Detection from
685 Video*. 2018.
- 686 74. Kim, H.-U. and C.-S. Kim, *CDT: Cooperative Detection and Tracking for Tracing Multiple Objects in Video
687 Sequences*, in *Computer Vision - Eccv 2016, Pt Vi*, B. Leibe, et al., Editors. 2016. p. 851-867.
- 688 75. Luo, H., et al., *Detect or Track: Towards Cost-Effective Video Object Detection/Tracking*. 2018.
- 689 76. Feichtenhofer, C., et al., *Detect to Track and Track to Detect*, in *2017 Ieee International Conference on
690 Computer Vision*. 2017. p. 3057-3065.
- 691 77. Sharma, V.K., B. Acharya, and K.K. Mahapatra, *Online Training of Discriminative Parameter for Object
692 Tracking-by-Detection in a Video*, in *Soft Computing in Data Analytics, Scda 2018*, J. Nayak, et al., Editors.
693 2019. p. 215-223.
- 694 78. Bertasius, G., L. Torresani, and J. Shi, *Object Detection in Video with Spatiotemporal Sampling Networks*.
695 2018.

- 696 79. Chen, K., et al., *Optimizing Video Object Detection via a Scale-Time Lattice*, in *2018 Ieee/Cof Conference on*
697 *Computer Vision and Pattern Recognition*. 2018. p. 7814-7823.
- 698 80. Wang, T., et al., *SCNN: A General Distribution Based Statistical Convolutional Neural Network with*
699 *Application to Video Object Detection*. Thirty-Third Aaai Conference on Artificial Intelligence / Thirty-
700 First Innovative Applications of Artificial Intelligence Conference / Ninth Aaai Symposium on
701 Educational Advances in Artificial Intelligence. 2019. 5321-5328.
- 702 81. Du, Y., et al., *Spatio-temporal self-organizing map deep network for dynamic object detection from videos*. 2017.
- 703 82. Xiao, F. and Y.J. Lee, *Video Object Detection with an Aligned Spatial-Temporal Memory*. Arxiv, 2017.
- 704 83. Jiang, Z., et al., *Video Object Detection with Locally-Weighted Deformable Neighbors*. Thirty-Third Aaai
705 Conference on Artificial Intelligence / Thirty-First Innovative Applications of Artificial Intelligence
706 Conference / Ninth Aaai Symposium on Educational Advances in Artificial Intelligence. 2019. 8529-
707 8536.
- 708 84. Zhu, H., et al., *Moving Object Detection With Deep CNNs*. Ieee Access, 2020. 8: p. 29729-29741.
- 709 85. Chin, T.-W., R. Ding, and D. Marculescu, *AdaScale: Towards Real-time Video Object Detection Using*
710 *Adaptive Scaling*. 2019.
- 711 86. Deng, J., et al., *ImageNet: A Large-Scale Hierarchical Image Database*, in *Cvpr: 2009 Ieee Conference on*
712 *Computer Vision and Pattern Recognition, Vols 1-4*. 2009. p. 248-255.
- 713 87. Lin, T.-Y., et al., *Microsoft COCO: Common Objects in Context*, in *Computer Vision - Eccv 2014, Pt V*, D.
714 Fleet, et al., Editors. 2014. p. 740-755.
- 715 88. Real, E., et al., *YouTube-BoundingBoxes: A Large High-Precision Human-Annotated Data Set for Object*
716 *Detection in Video*, in *30th Ieee Conference on Computer Vision and Pattern Recognition*. 2017. p. 7464-7473.
- 717 89. Damen, D.a.D., Hazel and Farinella, Giovanni and Fidler, Sanja and Furnari, Antonino and Kazakos,
718 Evangelos and Moltisanti, Davide and Munro, Jonathan and Perrett, Toby and Price, Will and et al., *The*
719 *EPIC-KITCHENS Dataset: Collection, Challenges and Baselines*. IEEE Transactions on Pattern Analysis and
720 Machine Intelligence, 2020: p. 1-1.
- 721 90. Perazzi, F., et al., *A Benchmark Dataset and Evaluation Methodology for Video Object Segmentation*, in *2016*
722 *Ieee Conference on Computer Vision and Pattern Recognition*. 2016. p. 724-732.
- 723 91. Wang, Y., et al., *CDnet 2014: An Expanded Change Detection Benchmark Dataset*, in *2014 Ieee Conference on*
724 *Computer Vision and Pattern Recognition Workshops*. 2014. p. 393-+.
- 725 92. Kristan, M., et al., *The Visual Object Tracking VOT2015 challenge results*. 2015 Ieee International
726 Conference on Computer Vision Workshop. 2015. 564-586.
- 727 93. Schindler, L.L.-T.a.A.M.a.I.R.a.S.R.a.K., *MOTChallenge 2015: Towards a Benchmark for Multi-Target*
728 *Tracking*. arXiv, 2015.
- 729 94. Han, G., et al., *Semi-Supervised DFF: Decoupling Detection and Feature Flow for Video Object Detectors*.
730 Proceedings of the 2018 Acm Multimedia Conference. 2018. 1811-1819.
- 731 95. Yang, Y., et al., *Semi-supervised Learning of Feature Hierarchies for Object Detection in a Video*, in *2013 Ieee*
732 *Conference on Computer Vision and Pattern Recognition*. 2013. p. 1650-1657.
- 733 96. Singh, K.K., et al., *Track and Transfer: Watching Videos to Simulate Strong Human Supervision for Weakly-*
734 *Supervised Object Detection*, in *2016 Ieee Conference on Computer Vision and Pattern Recognition*. 2016. p.
735 3548-3556.
- 736 97. Sharma, P., et al., *Unsupervised Incremental Learning for Improved Object Detection in a Video*, in *2012 Ieee*
737 *Conference on Computer Vision and Pattern Recognition*. 2012. p. 3298-3305.

- 738 98. Mao, H.a.Y., Xiaodong and Dally, Bill, *A Delay Metric for Video Object Detection: What Average Precision*
739 *Fails to Tell*. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019.
- 740 99. Dosovitskiy, A., et al., *FlowNet: Learning Optical Flow with Convolutional Networks*, in *2015 Ieee*
741 *International Conference on Computer Vision*. 2015. p. 2758-2766.
- 742 100. He, K., et al., *Deep Residual Learning for Image Recognition*, in *2016 Ieee Conference on Computer Vision and*
743 *Pattern Recognition*. 2016. p. 770-778.
- 744 101. Luo, C., et al., *Cosine Normalization: Using Cosine Similarity Instead of Dot Product in Neural Networks*.
745 2017.
- 746 102. J. Deng, Y.Z., B. Yu, Z. Chen, S. Zafeiriou, and a.D. Tao., *Speed/accuracy tradeoffs for object detection from*
747 *video*. 2017.
- 748 103. Hochreiter, S. and J. Schmidhuber, *Long short-term memory*. *Neural Computation*, 1997. 9(8): p. 1735-
749 1780.
- 750 104. Elkan, Z.C.L.a.J.B.a.C., *A Critical Review of Recurrent Neural Networks for Sequence Learning*. arXiv:
751 1506.00019, 2015.
- 752 105. Vaswani, A., et al., *Attention Is All You Need*, in *Advances in Neural Information Processing Systems 30*, I.
753 Guyon, et al., Editors. 2017.
- 754 106. Bahdanau, D., K. Cho, and Y. Bengio, *Neural Machine Translation by Jointly Learning to Align and Translate*.
755 *Computer Science*, 2014.
- 756 107. Xie, S., et al., *Aggregated Residual Transformations for Deep Neural Networks*, in *30th Ieee Conference on*
757 *Computer Vision and Pattern Recognition*. 2017. p. 5987-5995.
- 758 108. Chen, X., J. Yu, and Z. Wu, *Temporally Identity-Aware SSD With Attentional LSTM*. *Ieee Transactions on*
759 *Cybernetics*, 2020. 50(6): p. 2674-2686.
- 760 109. Chen, X., Z. Wu, and J. Yu, *TSSD: Temporal Single-Shot Detector Based on Attention and LSTM*, in *2018*
761 *Ieee/Rsj International Conference on Intelligent Robots and Systems*, A.A. Maciejewski, et al., Editors. 2018.
762 p. 5758-5763.
- 763 110. Zhu, H., et al., *Real-Time Moving Object Detection in High-Resolution Video Sensing*. *Sensors (Basel,*
764 *Switzerland)*, 2020. 20(12).
- 765



© 2020 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).