

Plug & Play Convolutional Regression Tracker for Video Object Detection

Ye Lyu
University of Twente

Michael Ying Yang
University of Twente

George Vosselman
University of Twente

Gui-Song Xia
Wuhan University

Abstract

Video object detection targets to simultaneously localize the bounding boxes of the objects and identify their classes in a given video. One challenge for video object detection is to consistently detect all objects across the whole video. As the appearance of objects may deteriorate in some frames, features or detections from the other frames are commonly used to enhance the prediction. In this paper, we propose a Plug & Play scale-adaptive convolutional regression tracker for the video object detection task, which could be easily and compatibly implanted into the current state-of-the-art detection networks. As the tracker reuses the features from the detector, it is a very light-weighted increment to the detection network. The whole network performs at the speed close to a standard object detector. With our new video object detection pipeline design, image object detectors can be easily turned into efficient video object detectors without modifying any parameters. The performance is evaluated on the large-scale ImageNet VID dataset. Our Plug & Play design improves mAP score for the image detector by around 5% with only little speed drop.

1. Introduction

Last several years have witnessed the fast development of the computer vision in scene understanding, especially the fundamental object detection task. Object detection is to simultaneously localize the bounding boxes of the objects and classify their class types in an image. Video object detection extends the task to videos. It requires the detector to utilize multiple frames in a video to provide consistent predictions over time, which is another fundamental task for computer vision. The large scale datasets and the convolutional neural networks have been two of the main propellers for such visual understanding tasks. Both of the tasks have received much attention since the introduction of ImageNet object detection (DET) challenge and ImageNet video object detection (VID) challenge in ILSVRC 2015 [35]. The video object detection is more challenging than the image object detection due to that the appearance of the objects

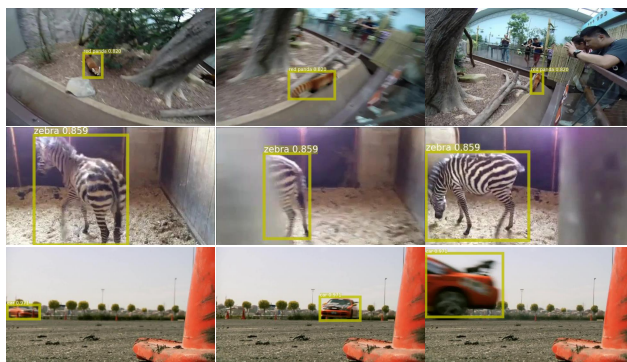


Figure 1. Examples of object detection in videos by our method. The scores are consistent over the long term even if the image quality decays or the objects are partially occluded or in rare poses.

may deteriorate significantly in some frames caused by motion blur, video defocus, part occlusion or rare poses. However, the rich context information in temporal domain provides the clues and the opportunity to classify the objects.

One basic video object detection method is to detect objects from individual image first and further apply a post-processing to link and re-score the object detection results in the video [15, 22]. The linking can be based on either the appearance similarity scores [45] or the external video object tracker [22]. The problem for such methods is that they require good enough initial image object detection results. If the initial detection is lost or the localization of the initial detection is not accurate, the post-processing performance may not be ideal. Feature aggregation is another widely used method for video object detection. By using features from multiple frames, a more temporally coherent and robust detection result can be achieved. The problem for short term feature aggregation [43, 1, 47, 41] is that they pre-define a limited range of frames for the detection in each frame and the long range coherency cannot be preserved and utilized. Long term feature aggregation achieves better results [42, 10, 37]. Feature aggregation normally cost more memory and time during inference as features from multiple frames are needed. Combining the detection and

the tracking is another direction worth attentions. Tracking could be an aid for the detection as it should propose better bounding boxes based on the detection results in the previous frames. As tracking looks for similarity between images, it is generally easier to track than to detect an object with deteriorated appearance between consecutive frames.

The tracker design for the video object detection task is different from it is for the video object tracking task. There are several techniques in object tracking task have to be rejected. As the detection network requires powerful recognition ability, the features extracted are often heavier than those for tracking. It would be preferable if the tracker could re-use the features from the detection. However, such deep feature may not be compatible with some state of the art object trackers, e.g. siamese region proposal networks [25, 24] and fully-convolutional siamese networks [2]. The performance would be heavily undermined by the different influence of paddings to the anchors in different spatial positions in a feature map. Another common practice in single-object tracking is to crop and resize the image patches of the template and the target to pre-defined sizes for valid network inputs. Such practice ensures that the features extracted are scale invariant to the real size of an object in the image. However, it is forbidden in the video object detection as it corrupts the features for object recognition.

To improve the combination of the detection and the tracking, we propose a novel tracker for the video object detection task, which is compatible with the detection networks. The **main contributions** of this paper are following:

- We have created an object tracker for video object detection task. It is very light-weighted, memory efficient, computationally efficient and compatible to deep features of the existing modern detection networks.
- Our tracker can be inserted into a well trained image detector in a Plug & Play style. Without harming the performance of the detector, the tracking functionality can be implanted into the model.
- Our new tracker performs in adaptive scales according to the size of the object being tracked, which makes it robust for tracking the object with large size variation in a video.
- We have designed a new video object detection pipeline to combine the advantages from both of the detection and the tracking. With better bounding box proposals and linkages through time, we improve the performance with better effectiveness and the efficiency.

2. Related work

Object detection in images is one of the fundamental tasks in computer vision. After the successful pioneer work [23] with deep neural networks in object detection, a number of one-stage and two-stage object detector networks have been proposed [13, 34, 4, 28, 27, 32, 33, 29, 7]. As an extension to object detection in images, object detection in videos has received much attention as well. Many methods have been proposed that introduce the idea of tracking, but very few achieve a real integration of detection and tracking networks. Instead, external trackers [22] or optical flows [48, 47, 41] or alternatives for tracking are required.

Seq-NMS [15] is a post-processing method for video object detection, which utilizes detection results from an object detector. Detections are linked through max score path finding under 0.5 IoU score constraint between boxes in consecutive frames. Linked detections are re-scored afterwards. Such constraint is not optimal as it may not hold for objects in fast movement.

T-CNN [22] proposes a deep learning framework that incorporates an external independent tracker [40] to link the detections, which makes the pipeline slow. In [21] tubelet proposal network is utilized to propose tubelet boxes for multiple frames simultaneously. Boxes in the same tubelet are linked. [45] learns an additional feature embedding to help link the detected objects to the corresponding tracklets. [6] uses a scale-time lattice to accelerate the speed for video object detection. The temporal propagation for detection is inferred from the motion history images [3].

D&T [12] brings the tracking into the detection. By utilizing the feature map correlations between frames, the model learns a box regression model from one image to another. The tracked pairs are used to boost the linking score between image detections. However, D&T [12] calculates the feature map correlations with a number of pre-defined position shifts and the feature map correlations have to be inferred for the whole image, which is very inefficient.

Another direction for video object detection is through feature aggregation. By gathering clues from multiple consecutive frames, semantic information from both spatial and temporal domain can be extracted simultaneously to boost the detection performance.

DFE [48], FGFA [47], MANet [41] adopt optical flows to warp the features for alignment. Modern deep learning based optical flow models, such as FlowNet [11], FlowNet2.0 [20] and LiteFlowNet [19] can process images in a very fast speed, which provide the chance for fast video object detector. However, optical flow models are normally trained with synthetic data and the performance of the tracking is limited by the domain discrepancy. STMN [43] aggregates the spatial-temporal memory for multiple frames according to the MatchTrans module, which is guided by the feature similarity between consecu-

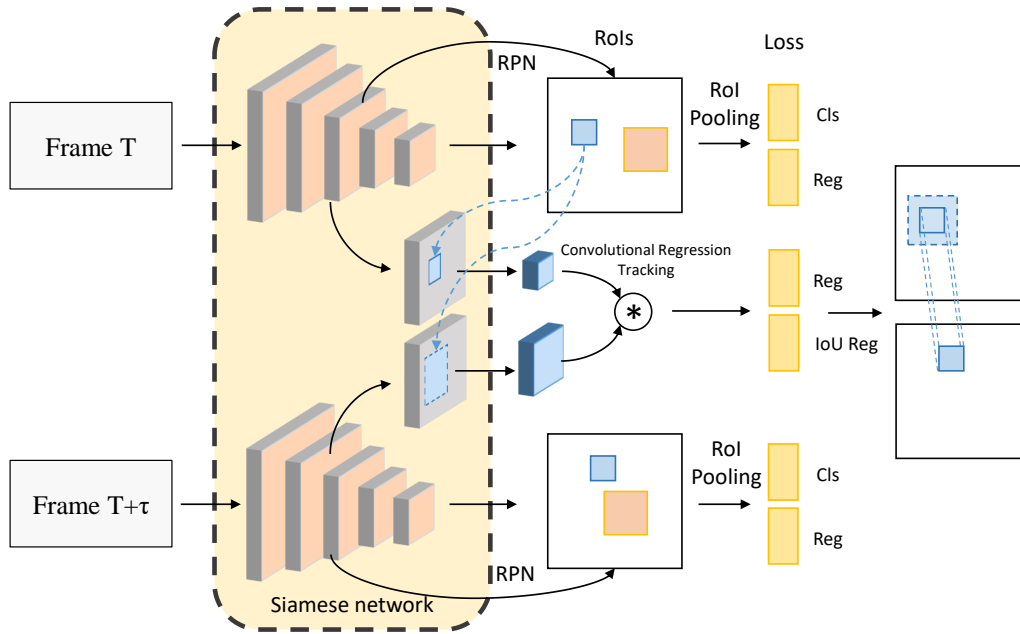


Figure 2. Architecture of our video object detection network. Our Plug&Play tracker reuses the features of the detection networks from both branches. The regional features within RoIs are pooled and sent to the tracker. The regional features from the two branches are convolved with each other for bounding box and IoU regression. (Details illustrated in Sec. 3).

tive frames. STSN [1] directly extracts the spatially aligned features by using deformable convolutions [8]. [14] adopts progressive sparse local attentions to propagate the features across frames, while [9] utilizes explicit external memory to accumulate information over time. The temporal connection of these methods may not be accurate as there lacks an explicit learning process for tracking. [37, 10, 42] are more powerful feature aggregation methods as they distil semantic information from longer sequences.

For tracking, siamese networks have received much attention recently. [2, 39] score the locations of objects by using feature correlation through a convolution operation between the template patch and the target patch. The idea is extended by [25, 24] with region proposal networks, which infer the object score and the box regression simultaneously for better box localization. GOTURN [18] adopts siamese network as feature extractor and uses fully connected layers to merge features for bounding box regression.

Our work is inspired by the ideas in D&T [12] and the siamese networks for tracking [18, 2, 39, 25].

3. Architecture Overview

In this section, we will introduce an overview of our model structure. The goal of our model is to plug the tracking network into the detection network without harming the performance of the image detector. The architecture design is shown in Fig. 2. Our model takes two

consecutive frames with a gap of τ (1 for testing) as inputs $I^t, I^{t+\tau} \in \mathbb{R}^{H \times W \times 3}$, followed by a siamese network for feature extraction. The two branches share the same weights to keep the identical feature extractors. To satisfy the need for object detection in complex scenes, we exploit powerful feature extraction backbones such as HRNet [38] and ResNeXt [17, 44]. The two models correspond to a light version detector and a heavy version detector. The extracted features from the siamese network are further sent to the detection branches and the tracking branch simultaneously. In a detection branch, regions of interest (RoIs) are proposed by the Region Proposal Network (RPN) [34], followed by the RoI Pooling layer [13] to extract the features within each RoI. The RoI-wise features are exploited for object classification and bounding box regression, which is the same as the faster-rcnn [34]. One branch of the siamese network plus the detection branch forms a standard two-stage object detector. In the tracking branch, the novel scale-adaptive convolutional regression tracker is utilized to predict the bounding box transformation from the first frame to the second frame. The tracker utilizes regional features from the two branches of the siamese network based on the RoIs to be tracked. During training, the RoIs are generated from the ground truth bounding boxes, while in the testing phase, the RoIs are the detected objects. Besides the bounding box regression from the first image to the second image, the tracker has another IoU score regression branch to estimate the quality of the bounding box regression.

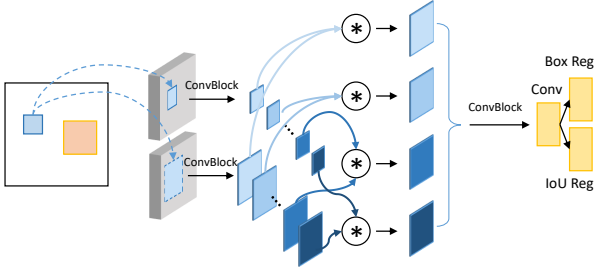


Figure 3. Illustration of depth-wise feature correlation in our tracker. The bounding box of an object determines the location where the tracker acquires the local features. The features of each channel from the first branch is convolved with the features of the same channel from the second branch. Convolutional blocks are inserted to adjust the features, each of which is comprised of a convolutional layer, a batch normalization layer and a relu layer.

4. Scale-adaptive convolutional regression tracker

Many trackers for the video object tracking task crop and resize the image patches to fixed sizes according to the sizes of the objects to be tracked [25, 24, 2]. Such standard sizes for the network inputs make feature extraction invariant to the sizes of the objects. However, image patch cropping and resizing are not applicable to the detection network as there maybe multiple objects of different sizes. Instead of regularizing the size of the network inputs, we aim to extract scale-adaptive features for tracking by reusing features from the shared backbone, which augments the detector in a Plug & Play style.

We extract regional features from both of the two branches based on the RoIs to be tracked. The RoI for the first branch marks the bounding box extent of an object. The width and height of the RoI bounding box are expanded k times for the second branch with the center point and the aspect ratio fixed, which marks the local area of interest to search for the object in the second frame. k is set to 3, indicating one object space to each side of the center object. RoIAlign [16] is adopted to pool the features from the two branches of the siamese network. To keep the same scale of the pooled features, the pooled feature size from the second branch is also k times the size of the first branch as shown in Fig. 4. The features pooled from outside the range of the image are set to zero. We adopt the backbone features from multiple stages for RoIAlign pooling. Features from different stages are resized to the same intermediate size with max pooling and interpolation operations as in [30]. Instead of averaging, we concatenate the resized features. For HRNet-W32 [38] backbone, all stages are adopted. For ResNeXt101 [44] backbone, features corresponding to middle 3 stages are applied.

D&T [12] utilizes feature correlation as clues for track-

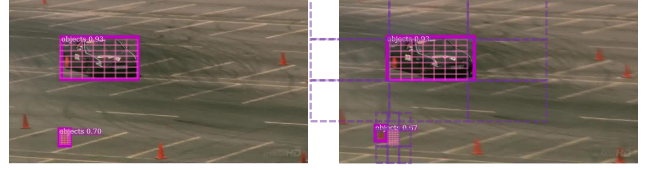


Figure 4. Scale-adaptive tracking feature extraction. Here shows a real example of two objects being tracked. The 7×7 features are pooled within the object bounding boxes from the first image. The 21×21 features are pooled within the search area from the second image, which is 3 times the size of the object.



Figure 5. Strategy for the bounding box selection. Bounding boxes can be inferred from both of the detection and the tracking process. We choose the tracked ones over the detected ones.

ing and encodes different translations into different channels. Our scale-adaptive tracker calculates feature correlation with a depth-wise convolution operation [24] between two feature patches from the two branches of the siamese network. Fig. 3 illustrates the process. Correlation of different translations are encoded into different spatial positions in our tracker. Convolutional blocks are inserted before and after the convolution for feature adjustment, each of which consists of a convolution layer, a batch normalization layer and a relu layer. The head of the tracker has a bounding box regression branch and an IoU score regression branch, which are two fully connected layers attached to a shared 2D convolution layer with 256 filters. Sigmoid function is applied to normalize the IoU score prediction. As we adopt a class-agnostic regression tracking, the output dimensions are 4 and 1 respectively for each object. Smooth L_1 loss is utilized for regression [13]. Our tracker learns to predict the target bounding boxes directly so that it can rectify the boxes during tracking.

The learning targets are defined by the bounding boxes to be tracked $b^t = (b_x^t, b_y^t, b_w^t, b_h^t)$ in time t , the predicted bounding boxes $p^{t+\tau} = (p_x^{t+\tau}, p_y^{t+\tau}, p_w^{t+\tau}, p_h^{t+\tau})$ in $t + \tau$ and the corresponding ground truth bounding boxes $g^{t+\tau} = (g_x^{t+\tau}, g_y^{t+\tau}, g_w^{t+\tau}, g_h^{t+\tau})$ in $t + \tau$. The target for bounding box regression $\Delta^{t+\tau} = (\Delta_x^{t+\tau}, \Delta_y^{t+\tau}, \Delta_w^{t+\tau}, \Delta_h^{t+\tau})$ is defined as,

$$\begin{aligned} \Delta_x^{t+\tau} &= \frac{g_x^{t+\tau} - b_x^t}{b_w^t}, \Delta_y^{t+\tau} = \frac{g_y^{t+\tau} - b_y^t}{b_h^t} \\ \Delta_w^{t+\tau} &= \ln \frac{g_w^{t+\tau}}{b_w^t}, \Delta_h^{t+\tau} = \ln \frac{g_h^{t+\tau}}{b_h^t} \end{aligned} \quad (1)$$

The target for IoU score regression is calculated by $b_{score}^{t+\tau} = IoU(p^{t+\tau}, g^{t+\tau})$, $b_{score}^{t+\tau} \in [0, 1]$. The predicted bounding boxes are inferred through predicted bounding box regression $\hat{\Delta}^{t+\tau} = (\hat{\Delta}_x^{t+\tau}, \hat{\Delta}_y^{t+\tau}, \hat{\Delta}_w^{t+\tau}, \hat{\Delta}_h^{t+\tau})$. $p^{t+\tau}$ can be calculated as,

$$\begin{aligned} p_x^{t+\tau} &= \hat{\Delta}_x^{t+\tau} b_w^t + b_x^t, p_y^{t+\tau} = \hat{\Delta}_y^{t+\tau} b_h^t + b_y^t \\ p_w^{t+\tau} &= \exp(\hat{\Delta}_w^{t+\tau}) b_w^t, p_h^{t+\tau} = \exp(\hat{\Delta}_h^{t+\tau}) b_h^t \end{aligned} \quad (2)$$

Our feature extraction design for tracking has several advantages. First, only resized local features are utilized, which makes the tracker memory efficient. Second, features extracted from the two branches are of the same scale and are in scale to the size of an object. The feature correlation can be calculated in an adaptive scale and range, which aligns with the scale-adaptive learning target for bounding box regression. Finally, the tracker is light-weighted as it reuses the features from the backbone network.

5. Combination of detection and tracking

Our model has the functionality for both of the detection and the tracking. In this section, we introduce how the detection and the tracking are combined to solve the video object detection task. The aim of detection is to find the newly appeared objects in an image, while the tracking is to better localize the objects across the frames.

Detect to track. The two-stage detection network has two object proposal stages, the RoI proposals by RPN and the object detections by R-CNN. We base our tracking on top of the object detections instead of the RoI proposals as the R-CNN provides much cleaner and more accurate objects for tracking. The design of RPN is to provide proposals with high enough recall rate and it is not time or memory efficient to track hundreds of RoIs. Setting a score threshold for foreground proposal selection would not be good either as the scores and bounding boxes would not be accurate enough due to the anchor based design. The detection results from R-CNN is more robust. If an object is identified in an image with a high enough class score (0.03 in our case), it would be a candidate for tracking in the next image.

Track to detect. Tracking could aid the detection by providing better object boxes. We first filter the tracked objects. If the predicted tracking IoU score is too small for an object (0.5 in our case), the object is discarded. As multiple objects may overlap with each other during tracking, we apply a non-maximum suppression (0.7 IoU threshold in our case) to the tracked boxes based on the IoU score. The selected objects are used in the video object detection. A tracking first detection (TFD) strategy is applied. When an object is identified by both of the detection and the tracking, we choose the tracked one over the detected one as shown in Fig. 5. The non-maximum suppression in detection favours the objects with higher class scores instead of the objects

with more accurate bounding boxes. We favour the better bounding boxes by adopting the TFD strategy to help acquire better object linking across frames. Only the newly detected objects that have IoUs with the existing tracked ones lower than a threshold T_{merge}^{mms} are reserved. During the inference across all the image frames in a video, the tracked boxes are saved. Detections are further refined by averaging re-scoring and non-maximum suppression as in [15].

6. Experiments

6.1. Dataset and evaluation

Our method is evaluated on the ImageNet [35] object detection from video (VID) dataset. There are 3862 training and 555 validation videos with objects from 30 classes labelled for the task. The ground truth annotations contain the bounding box, the class ID and the track ID for each object. The performance of the algorithm is measured with mean average precision (mAP) score over the 30 classes on the validation set as it is in [47, 41, 12, 22, 21, 43, 1]. In addition to the ImageNet VID dataset, the ImageNet object detection (DET) dataset has 200 classes, which include all the 30 classes in the ImageNet VID dataset as well. We follow the common practise by utilizing the data from both of the DET and the VID dataset [47, 41, 12, 22, 21, 43, 1].

6.2. Configurations

Image training. In the first stage, we train the detection parts of our video object detector in the same way as a standard object detector. The training samples are from both of the DET and the VID dataset. To balance the classes in the DET dataset, we sample at most $2K$ images from each of the 30 categories to get our DET image set (53K images). To balance the VID videos, which have large sequence length variations, we evenly sample 15 frames from each of the video sequence to get our VID image set (57K images). The combined DET+VID image set is used for detector training. We apply SGD optimizer with a learning rate of 10^{-3} for the first 90K iterations and 10^{-4} for the last 45K iterations. The training batch is set to 8 images that are distributed among 4 gpus. In both of the training and the testing, we apply a single scale with the shorter dimension of the images to be 600 pixels. During training, only random left-to-right flip is used for data augmentation ¹.

Video training. In this stage, we further train our tracking parts with the image pairs from the ImageNet VID dataset. We randomly select two consecutive images with a random temporal gap from 1 to 9 frames. As there is no causal reasoning involved, we randomly reverse the sequence order to gain more variety of translations. The RoIs for tracking $R = (R_x, R_y, R_w, R_h)$ are generated by resizing and shifting the ground truth bounding boxes

¹Please check the supplementary material for detailed configurations.

$g = (g_x, g_y, g_w, g_h)$ as in Eq. 3. The coefficients $\delta = (\delta_x, \delta_y, \delta_w, \delta_h)$ are sampled from uniform distributions U . $\delta_x, \delta_y \in U[-1.0, 1.0]$ and $\delta_w, \delta_h \in U[0.5, 1.5]$. For each ground truth object, we sample 256 RoIs and randomly select 128 RoIs from those satisfying the constraint of $IoU(R, g) < 0.5$.

$$\begin{aligned} R_x &= \delta_x g_w + g_x, & R_y &= \delta_y g_h + g_y, \\ R_w &= \delta_w g_w, & R_h &= \delta_h g_h, \end{aligned} \quad (3)$$

We freeze the backbone parts to train the tracking parts only in order to retain the accuracy for detection. RPN and R-CNN are not included either. The model is trained with SGD optimizer with a learning rate of 10^{-3} for the first $80K$ iterations and 10^{-4} for the next $40K$ iterations. We apply a batch of 16 image pairs for training that are distributed among 4 gpus. The images are resized to the same single scale as the image training step.

Testing. In the testing stage, we select the detected objects with the class scores higher than 0.03 and apply an IoU threshold of 0.45 for the final detection output. The single scale testing is utilized with the shorter side of images resized to 600 pixels.

Implementations. Our model is implemented with pytorch [31] and integrated with MMDetection [5].

6.3. Results

We compare several major competitive video object detection algorithms in Tab. 1. FGFA [47] and MANet [41] utilize optical flow to guide linking between frames, but they cannot achieve a good balance between the mAP score and the speed. STMN [43] and STSN [1] aggregate information from multiple frames, which limits the speed greatly. By bringing a very light-weighted tracker into the model, our light and heavy models achieve scores comparable to some of the state-of-the-art methods.

6.4. Ablation study

To test the effectiveness of bringing the tracking into the model, we perform ablation study by gradually adding the components into the model. We start with the per-image detection model. The faster-RCNN with HRNet backbone is adopted as the basic model. We first test the standard detection result without any aid from the tracking. The performance score is shown in Tab. 6. The per-image detection achieves a decent mAP score of 73.2%. We further add Seq-NMS [15] to see the performance of linking and re-scoring based solely on the detection results. As in [15], we set the IoU threshold for boxes linking constraint to be 0.5 and IoU threshold for detection NMS to be 0.45. The average precision scores improve for all the categories and the mAP score has increased by 2.3%. We further add tracking into our model by adopting our TFD video object detection

Methods	Temporal link	mAP (%)	FPS
FGFA [47]	optical flow	78.4	1.15
FGFA+ [47]	optical flow	80.1	1.05
MANet [41]	optical flow	78.1	4.96
MANet+ [41]	optical flow	80.3	-
STMN [43]	STMM	80.5	1.2
STSN [1]	DCN	78.9	-
STSN+ [1]	DCN	80.4	-
D&T [12]	box regression	79.8	7.09
PSLA+ [14]	attention	81.4	5.13
EMN [9]	memory	79.3	8.9
EMN+ [9]	memory	81.6	-
Ours(HRNet-w32)	box regression	78.6	11
Ours(ResNeXt101*)	box regression	81.1	5.6

Table 1. Comparisons among different video object detection methods. + stands for Seq-NMS [15]. - stands for not provided. * means with FPN [26] and DCN [8].

Methods	mAP (%)
HRNet-w32	73.2
HRNet-w32+Seq-NMS	75.5
HRNet-w32+TFD(0.3)+Seq-NMS	77.8
HRNet-w32+TFD(0.7)+Seq-NMS	78.6
ResNeXt101*	76.2
ResNeXt101*+TFD(0.7)+Seq-NMS	80.4
ResNeXt101*+TFD(0.7)+Seq-Track-NMS	81.1

Table 2. Ablation study of our method. TFD stands for tracking first detection. * stands for with FPN [26] and DCN [8].

strategy. During the training of the tracking modules, we freeze the parameters of the feature extraction backbone, the RPN and the R-CNN in order to control the performance of the object detector. We compare two values for merging NMS IoU threshold T_{merge}^{nms} , 0.3 and 0.7 (marked as TFD(0.3) and TFD(0.7) in Tab. 6). The mAP scores have increased another 2.3% and 3.1% respectively. The mAP score is higher with $T_{merge}^{nms} = 0.7$, showing that the video object detector still benefits more from the denser object proposals. The TFD strategy is very effective to help improve the quality of linking and re-scoring. By now, we have improved the performance of the object detector by a large margin (+5.4% mAP) without modifying any parameter of an object detector. With heavier ResNeXt101 [44] backbone plus FPN [26] and DCN [8], the TFD+Seq-NMS improves the detector by a large margin (+4.2% mAP) from 76.2% to 80.4%. Seq-NMS links boxes across frames under constraint $IoU(b^t, b^{t+1}) > 0.5$, which fails if there is large position translation between consecutive frames. We improve the constraint to be $IoU(p^{t+1}, b^{t+1}) > 0.5$, where p^{t+1} is the predicted box from b^t by the tracker. The improved re-scoring method (Seq-Track-NMS) provides another 0.7% mAP score boost.

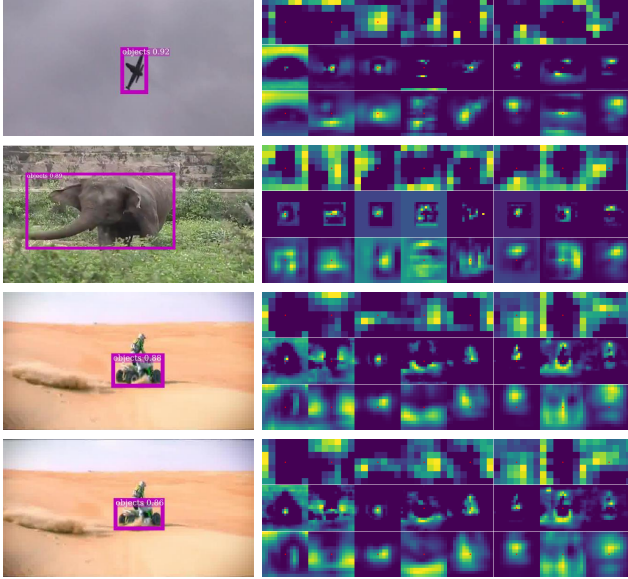


Figure 6. Correlation feature visualization. Images on the left show the objects being tracked. Feature maps on the right are the correlation features and the output features. For each image, the upper two rows are the correlation features from the two branches while the bottom row shows the convolution output. The randomly selected 8 channels are ordered in 8 columns. The center of a feature map is marked with a red dot for spatial reference.

6.5. Visualization of the correlation features

To make sure that the tracker utilizes the features from both of the branches rather than being dominated by one branch, we randomly sample the same 8 channels from the features before and after the correlation operation for visualization to check how the features for tracking are produced. The features are shown in Fig. 6. The examples are air-plane, elephant and motorcycles. The two motorcycle examples are two consecutive images, where the second image shows an inverse sequence order for tracking. It should be noted that the template and the target patches are encoded in the same scale but of different sizes. It is interesting to notice that the center of mass of the template features are shifted to all different positions. The features of different objects are distinctive and the correlation output is greatly affected from both of the branches. By examining the examples of the two motorcycles, which have opposite moving patterns, it could be seen that the features from the target branch determine the tracking prediction while the features from the template branch are more similar. In conclusion, the tracking prediction is mainly affected by the relative position for an object. The template features of different objects would encode the target features differently. The features from both of the branches affect the tracking prediction in the same time.

Components	Backbone	RPN	R-CNN	Tracker
Time (ms)	43	10	7	3

Table 3. Time cost comparison of different components.

Model	Detector	Detector+Tracker
Memory (GB)	1.59	1.65

Table 4. Memory cost comparison of different models.

Backbone	Detector	TFD	Seq-NMS	FPS
HRNet-w32	✓	✗	✗	15
HRNet-w32	✓	✓	✗	12
HRNet-w32	✓	✓	✓	11
ResNeXt101*	✓	✗	✗	6.9
ResNeXt101*	✓	✓	✓	5.6

Table 5. Time efficiency comparison of the pipelines. * stands for with FPN [26] and DCN [8].

6.6. The efficiency of time and memory

In this section, the efficiency of time and memory are examined. We take HRNet-w32 as an example for the following analysis. All the experiments are conducted with a single Titan X (Pascal) GPU during testing stage. The design of our tracking module is very light-weighted in both of the time and the memory. We first test the time efficiency of different components in our model. The approximate time costs are reported in Tab. 3. Our additional tracker is very efficient and costs only an extra time of 3 ms, which is lighter than the RPN or the R-CNN.

We further examine the running GPU memory cost of our model. We compare the GPU memory consumption with and without the tracker as shown in Tab. 4. Our tracker requires only another 60 MB to run, which is very memory efficient.

For the video object detection, we examine the running speed of our pipeline and analyse the effect of different components. The speed is measured in frames per second (FPS). The detector with HRNet-w32 backbone runs at 15 FPS and our tracker embedded pipeline runs at 12 FPS. The additional Seq-NMS slows down our pipeline slightly to 11 FPS. As our detector with ResNeXt101 backbone is combined with FPN [26] and DCN [8], it runs relatively slower but still faster than methods like FGFA+ [47], STMN [43] and MANet+ [41].

6.7. Qualitative results

Our tracker learns to track objects in different circumstances. The tracking examples of motion blur, rare pose, partial occlusion and multiple objects are shown in Fig. 7. By incorporating tracking into the detection, our video object detector can achieve very long term detection consis-



Figure 7. Tracking examples by our Plug & Play tracker. Our tracker can track single or multiple objects and can handle problems like motion blur, partial occlusion and rare object poses. The objects are labelled with IoU regression scores.

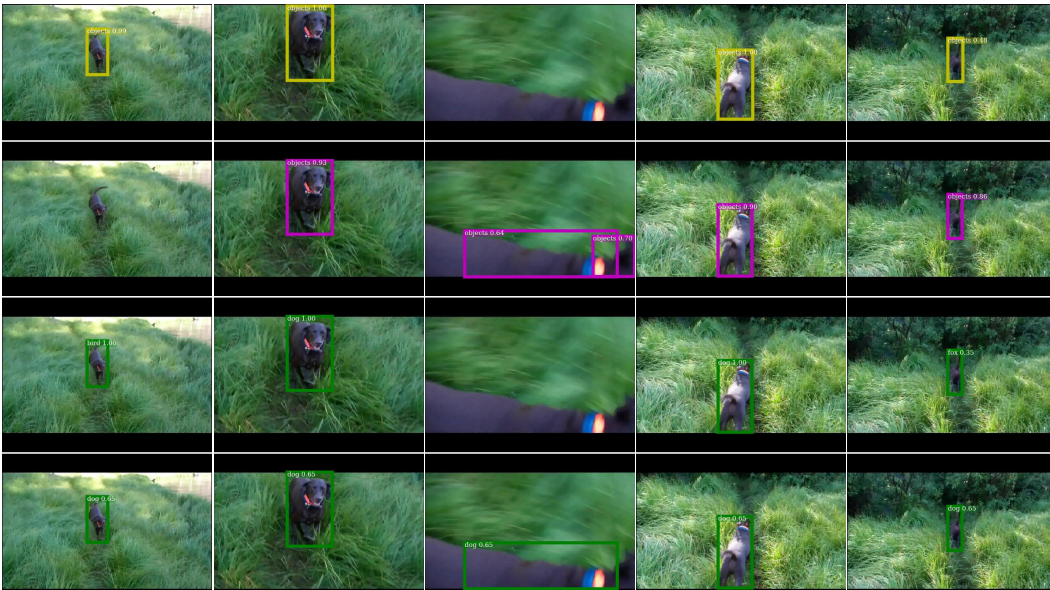


Figure 8. A video object detection example. The 0_{th} , 20_{th} , 40_{th} , 60_{th} , and 80_{th} images are shown. The first row shows the proposed object from the detector. The second row shows the proposed object from the tracker. The third row shows the detection from the combined proposals. The fourth row shows the re-scored detection results. Only the objects with scores above 0.2 are shown.

tency across frames. Fig. 8 shows an example of a running dog in a long sequence. This example clearly shows how the detector and the tracker collaborate. The detector finds new objects while the tracker follows the objects. The tracker provides better boxes for linking and re-scoring. Without re-scoring, the detection could be wrong or weak as shown in the third row (The dog is wrongly classified as a bird in the first column or a fox in the last column). After re-scoring, the long term detection consistency can be achieved as shown in the fourth row².

²More qualitative results are shown in the supplementary material.

7. Conclusion

We have proposed a Plug & Play convolutional regression tracker that augments the image detectors for the object detection in videos. The tracker makes use of the deep features from the image detectors for tracking with very little extra memory and time cost. The light-weighted tracker can track a single object or multiple objects, and handle problems e.g. motion blur, defocus, rare poses or partial occlusions. With our tracking first detection strategy and the improved Seq-NMS [15] linking and re-scoring method, the performance of our detector improves by a large margin.

References

- [1] Gedas Bertasius, Lorenzo Torresani, and Jianbo Shi. Object detection in video with spatiotemporal sampling networks. In *ECCV*, September 2018. 1, 2, 5, 6
- [2] Luca Bertinetto, Jack Valmadre, João F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. 2016. 2, 3, 4
- [3] A. F. Bobick and J. W. Davis. The recognition of human movement using temporal templates. *TPAMI*, 23(3):257–267, 2001. 2
- [4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *CVPR*, pages 6154–6162, 2018. 2
- [5] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 6
- [6] Kai Chen, Jiaqi Wang, Shuo Yang, Xingcheng Zhang, Yuanjun Xiong, Chen Change Loy, and Dahua Lin. Optimizing video object detection via a scale-time lattice. In *CVPR*, June 2018. 2
- [7] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *NeurIPS*. 2016. 2
- [8] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, Oct 2017. 3, 6, 7, 10, 11
- [9] Hanming Deng, Yang Hua, Tao Song, Zongpu Zhang, Zhenhui Xue, Ruhui Ma, Neil Robertson, and Haibing Guan. Object guided external memory network for video object detection. In *ICCV*, October 2019. 3, 6
- [10] Jiajun Deng, Yingwei Pan, Ting Yao, Wengang Zhou, Houqiang Li, and Tao Mei. Relation distillation networks for video object detection. In *ICCV*, October 2019. 1, 3
- [11] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, December 2015. 2
- [12] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Detect to track and track to detect. In *ICCV*, Oct 2017. 2, 3, 4, 5, 6
- [13] Ross Girshick. Fast r-cnn. In *ICCV*, pages 1440–1448, 2015. 2, 3, 4
- [14] Chaoxu Guo, Bin Fan, Jie Gu, Qian Zhang, Shiming Xiang, Veronique Prinet, and Chunhong Pan. Progressive sparse local attention for video object detection. In *ICCV*, October 2019. 3, 6
- [15] Wei Han, Pooya Khorrami, Tom Le Paine, Prajit Ramachandran, Mohammad Babaeizadeh, Honghui Shi, Jianan Li, Shuicheng Yan, and Thomas S. Huang. Seq-nms for video object detection. *CoRR*, abs/1602.08465, 2016. 1, 2, 5, 6, 9
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *ICCV*, Oct 2017. 4
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, June 2016. 3
- [18] David Held, Sebastian Thrun, and Silvio Savarese. Learning to track at 100 fps with deep regression networks. In *ECCV*, 2016. 3
- [19] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite-flownet: A lightweight convolutional neural network for optical flow estimation. In *CVPR*, June 2018. 2
- [20] Eddy Ilg, Nikolaus Mayer, Tomoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, July 2017. 2
- [21] Kai Kang, Hongsheng Li, Tong Xiao, Wanli Ouyang, Junjie Yan, Xihui Liu, and Xiaogang Wang. Object detection in videos with tubelet proposal networks. In *CVPR*, July 2017. 2, 5
- [22] Kai Kang, Hongsheng Li, Junjie Yan, Xingyu Zeng, Bin Yang, Tong Xiao, Cong Zhang, Zhe Wang, Ruohui Wang, Xiaogang Wang, and Wanli Ouyang. T-cnn: Tubelets with convolutional neural networks for object detection from videos. *IEEE Transactions on Circuits and Systems for Video Technology*, 28:2896–2907, 2018. 1, 2, 5
- [23] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1097–1105, 2012. 2
- [24] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. Siamrpn++: Evolution of siamese visual tracking with very deep networks. In *CVPR*, June 2019. 2, 3, 4
- [25] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with siamese region proposal network. In *CVPR*, June 2018. 2, 3, 4
- [26] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, July 2017. 6, 7, 10, 11
- [27] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017. 2
- [28] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *CVPR*, pages 8759–8768, 2018. 2
- [29] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016. 2
- [30] Jiangmiao Pang, Kai Chen, Jianping Shi, Huajun Feng, Wanli Ouyang, and Dahua Lin. Libra r-cnn: Towards balanced learning for object detection. In *CVPR*, June 2019. 4
- [31] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NeurIPS Autodiff Workshop*, 2017. 6

- [32] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, pages 779–788, 2016. 2
- [33] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *CVPR*, pages 7263–7271, 2017. 2
- [34] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, pages 91–99, 2015. 2, 3
- [35] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 1, 5
- [36] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, June 2016. 10
- [37] Mykhailo Shvets, Wei Liu, and Alexander C. Berg. Leveraging long-range temporal relationships between proposals for video object detection. In *ICCV*, October 2019. 1, 3
- [38] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, June 2019. 3, 4, 10
- [39] Jack Valmadre, Luca Bertinetto, Joao Henriques, Andrea Vedaldi, and Philip H. S. Torr. End-to-end representation learning for correlation filter based tracking. In *CVPR*, July 2017. 3
- [40] Lijun Wang, Wanli Ouyang, Xiaogang Wang, and Huchuan Lu. Visual tracking with fully convolutional networks. In *ICCV*, December 2015. 2
- [41] Shiyao Wang, Yucong Zhou, Junjie Yan, and Zhidong Deng. Fully motion-aware network for video object detection. In *ECCV*, September 2018. 1, 2, 5, 6, 7
- [42] Haiping Wu, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Sequence level semantics aggregation for video object detection. In *ICCV*, October 2019. 1, 3
- [43] Fanyi Xiao and Yong Jae Lee. Video object detection with an aligned spatial-temporal memory. In *ECCV*, September 2018. 1, 2, 5, 6, 7
- [44] Saining Xie, Ross Girshick, Piotr Dollar, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, July 2017. 3, 4, 6, 10
- [45] Zheng Zhang, Dazhi Cheng, Xizhou Zhu, Stephen Lin, and Jifeng Dai. Integrated object detection and tracking with tracklet-conditioned detection. *CoRR*, abs/1811.11167, 2018. 1, 2
- [46] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. *CoRR*, abs/1811.11168, 2018. 10
- [47] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. Flow-guided feature aggregation for video object detection. In *ICCV*, Oct 2017. 1, 2, 5, 6, 7
- [48] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. Deep feature flow for video recognition. In *CVPR*, July 2017. 2

Appendix

Additional configuration details.

Features extraction for tracking: The backbone features from the light HRNet-w32 [38] model and the FPN [26] features from the heavy ResNeXt101 [44] model are utilized as input for tracker. Features are resized and concatenated. Our ResNeXt101 model has the cardinality of 32 and the base width of 4.

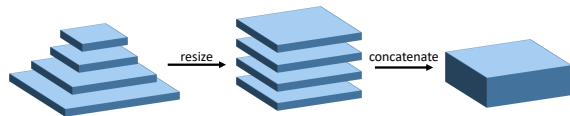


Figure 9. The feature from the HRNet-w32 backbone for tracking. The features are spatially resized to the size of the features in the second stage.

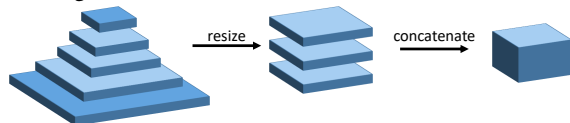


Figure 10. The features from the FPN (coupled with ResNeXt101) for tracking. The features from the middle 3 stages are resized to the size of the features in the third stage.

Region Proposal Network (RPN): The anchors in RPN have 3 aspect ratios (0.5, 1, 2) spanning 5 scales (32, 64, 128, 256, 512). For the ResNeXt model with FPN, 5 scales are distributed to the 5 stages in the FPN pyramid.

RoIAlign Pooling: The RoI pooling module has a size of 7×7 . For the tracker, the pooled feature is the average of the features in a bin. For the detector, the pooled feature is the average of the nearest 4 features.

R-CNN: The R-CNN has a bounding box regression branch and a logistic regression branch for classification. The two branches have 2 shared fully connected layers with 1024 filters, attached with 1 fully connected layer in each branch for their own purpose.

Deformable-ConvNets (DCN): The DCN [8] with 32 groups and modulation [46] is applied for stage 3,4,5 in ResNeXt [44] model.

Image training: Online hard example mining (OHEM) [36] is utilized for R-CNN training.

Detailed results

We additionally provide detailed mAP scores for ablation study as shown in Table 6.

More qualitative examples

We provide more qualitative video object detection examples in Fig. 11, Fig. 12 and Fig. 13 to show the effectiveness of our method. The tracker may track background objects as well, which is good as their scores for foreground objects will be smoothed to be lower after re-scoring.

Methods	airplane	antelope	bear	bicycle	bird	bus	car	cattle	dog	d. cat	elephant	fox	g. panda	hamster	horse	lion
HRNet-w32	83.7	82.8	79.6	73.4	72.0	68.3	58.0	68.5	65.5	73.4	75.8	86.1	81.4	91.9	69.2	48.2
HRNet-w32+Seq-NMS	83.8	85.1	83.8	73.8	72.5	70.0	59.0	70.8	67.9	78.5	76.7	88.7	81.8	96.2	70.6	58.3
HRNet-w32+TFD(0.3)+Seq-NMS	80.9	84.6	86.9	73.6	74.0	74.0	56.4	82.6	72.0	87.8	76.0	96.1	82.3	98.0	74.3	62.3
HRNet-w32+TFD(0.7)+Seq-NMS	87.0	86.0	85.6	76.6	71.7	75.0	56.6	80.8	72.5	89.8	80.4	95.4	82.4	98.8	79.3	63.8
ResNeXt101*	87.9	78.7	77.5	73.3	71.8	82.5	60.5	73.6	70.7	82.0	75.6	90.7	86.0	91.6	74.4	57.3
ResNeXt101*+TFD(0.7)+Seq-NMS	88.5	84.3	81.8	74.6	72.5	89.0	58.7	85.2	79.7	92.3	78.6	98.7	86.6	98.8	80.1	66.7
ResNeXt101*+TFD(0.7)+Seq-Track-NMS	88.3	84.7	83.2	74.6	72.8	88.8	58.5	85.3	80.5	92.3	78.7	98.7	86.2	98.9	80.4	71.2

Methods	lizard	monkey	motorcycle	rabbit	red panda	sheep	snake	squirrel	tiger	train	turtle	watercraft	whale	zebra	mAP (%)
HRNet-w32	82.4	47.1	81.2	70.8	81.2	55.1	73.5	56.2	89.7	78.0	79.2	63.8	70.1	91.0	73.2
HRNet-w32+Seq-NMS	84.0	49.4	83.3	75.0	87.9	57.3	74.2	57.2	90.2	78.3	80.7	65.2	72.5	91.0	75.5
HRNet-w32+TFD(0.3)+Seq-NMS	86.8	48.6	83.9	80.5	94.8	55.3	74.9	63.1	90.0	80.3	82.7	70.3	67.6	93.6	77.8
HRNet-w32+TFD(0.7)+Seq-NMS	88.3	51.7	87.8	80.2	93.2	58.3	73.7	57.3	89.9	82.3	83.0	72.2	67.1	91.0	78.6
ResNeXt101*	79.2	52.2	82.2	74.9	71.5	61.8	79.7	58.2	91.9	86.0	81.2	67.4	73.6	91.5	76.2
ResNeXt101*+TFD(0.7)+Seq-NMS	83.0	55.1	87.7	82.0	76.4	64.4	77.1	69.2	91.7	86.4	85.2	70.5	72.9	94.0	80.4
ResNeXt101*+TFD(0.7)+Seq-Track-NMS	82.8	55.2	87.7	82.1	90.8	64.4	76.9	69.1	91.7	86.4	85.3	70.7	72.4	94.6	81.1

Table 6. Ablation study of our method. TFD stands for tracking first detection. * stands for with FPN [26] and DCN [8].

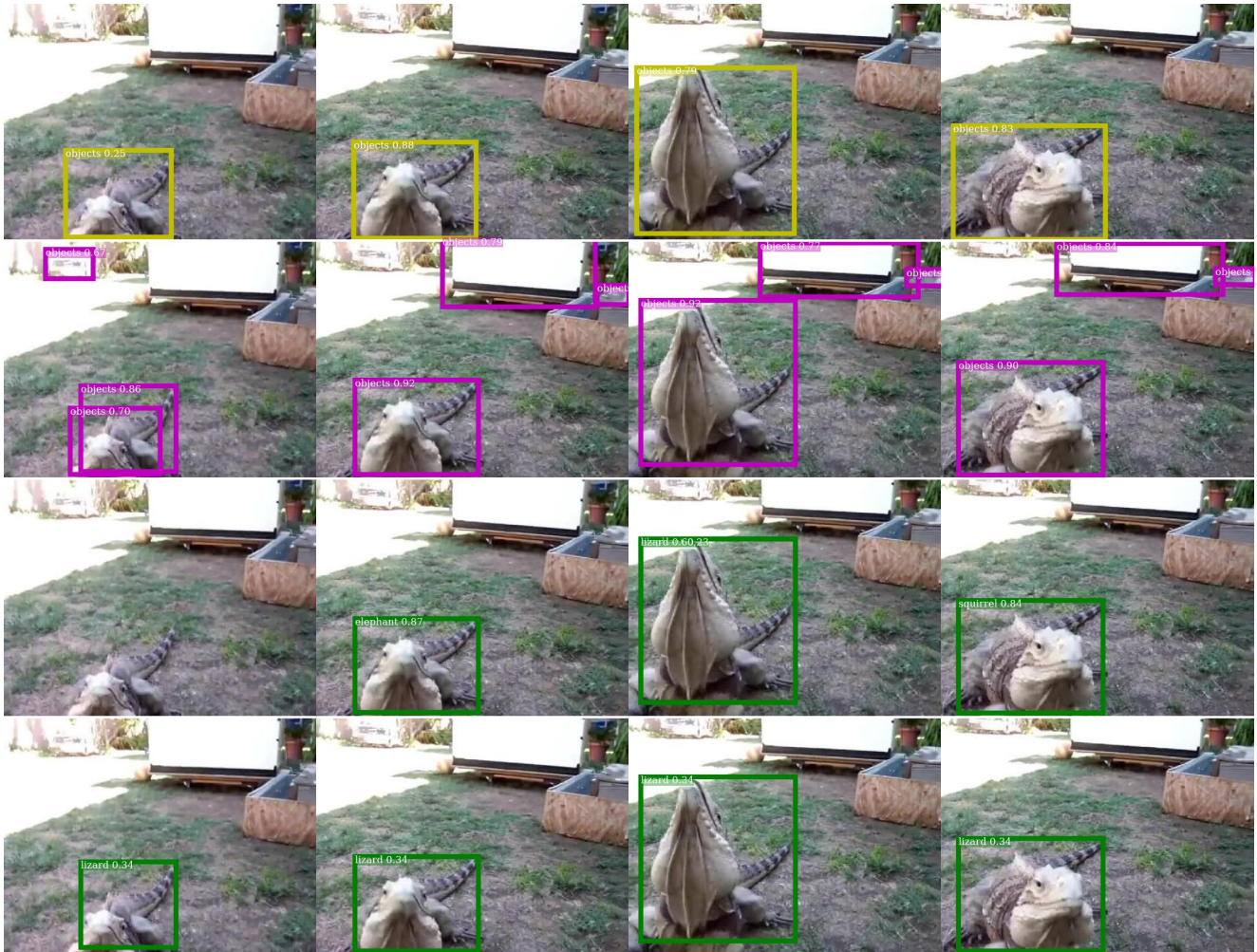


Figure 11. Example of lizard. The first row shows the proposed object from the detector. The second row shows the proposed object from the tracker. The third row shows the detection from the combined proposals. The fourth row shows the re-scored detection results. Only the objects with scores above 0.2 are shown. The weak or wrong detections shown in the third row are rectified as shown in the fourth row.



Figure 12. Example of red panda. The first row shows the proposed object from the detector. The second row shows the proposed object from the tracker. The third row shows the detection from the combined proposals. The fourth row shows the re-scored detection results. Only the objects with scores above 0.2 are shown.



Figure 13. Example of bus and cars. The first row shows the proposed object from the detector. The second row shows the proposed object from the tracker. The third row shows the detection from the combined proposals. The fourth row shows the re-scored detection results. Only the objects with scores above 0.2 are shown.