# 3D Object Detection and tracking in Point Cloud Flow

BMVC 2019 Submission # ??

#### Abstract

Recent approaches for 3D object detection has made great progresses due to the evolution in deep learning. However, previous works are mostly based on single frame point cloud or image, information between point cloud frames is almost not utilized. In this paper, we try to leverage the temporal information in point cloud flow and explore 3D object detection and tracking based on 3D data flow. Towards this goal, we set up a ConvNet architecture that can associate multi-frame image and LiDAR data to produce accurate 3D detection boxes and trajectories. Notably, a correlation module is introduced to capture object co-occurrences across time, and a multi-task objective for frame-based object detection and across-frame track regression is used, therefore, the network can performs detection and tracking simultaneously. Our proposed architecture is shown to produce competitive results on the KITTI Object tracking datasets. Code and models will be available soon.

#### 1 Introduction

A lot of attention has been paid to 3D object detection because of the rapid development of autonomous driving industry in recent years. Although object detection in images has made tremendous progress due to the emergence of deep learning [13, 14], 15] and their region based descendants [16, 14], extending 2D approaches to 3D scene is extremely hard. This is mainly because of the curse of dimensionality and the sparsity of 3D data.

In spite of the difficulty in 3D object detection, many works have been carried out. Recent approaches are usually done in three fronts: image based, point cloud based and fusion of image and point cloud. These methods have achieved remarkable performance but are limited to single frame input, applying existing detection networks on individual frames will loss the consistency and difference between frame and introduce unaffordable computational cost for most applications.

Compare to single point cloud frame, point cloud flow data is more natural and straightforward for most situation. Thus Fast and accurate point cloud flow object detection is crucial for autonomous driving. Similar to extend 2D object detection methods to 3D situation, we also can extend video object detection approaches to 3D point cloud flow object detection.

Most modern computer vision approaches to video object detection require flow estimation, which is a fundamental task in video analysis. For example, a series work in [26, 27] associate vision feature and optical flow to build an accurate and end-to-end learning framework for video object detection. However, applying video object detection framework to autonomous driving is hard, because video object detection suffers from motion blur and

<sup>© 2019.</sup> The copyright of this document resides with its authors. It may be distributed unchanged freely in print or electronic forms.

partial occlusion, which are conventional in driving scenarios. While another choice is to utilize LiDAR data, since point cloud provides an accurate spatial representation of the world 047 allowing for precise positioning of objects of interest, motion blur and occlusion problem 048 can be avoided in 3D object detection. However, LiDAR does not capture appearance well 049 when compared to the richness of images, moreover, accurate 3D scene flow estimation from 050 point clouds is tremendously tough. Although some approaches such as [III, III] have been 051 presented to learn 3D motion field in the world, it is hard to implementation in large scenes. 052 Thus the challenge is, can we propose an approach that accurately do 3D object detection in 053 point cloud flow without 3D scene flow?

Inspired by [1], we transform AVOD [15] structure into a two stream network embedding 055 with a correlation module, named Bi-AVOD, which takes two adjacent key frames as input 056 and predicts location and orientation of object as well as their local displacement. Noting 057 that Bi-AVOD is an aggregate view object detection architecture capable of fusing different 058 features in image and point cloud, thus its input includes two adjacent images in front view 059 and two adjacent BEV (bird eye view) from LIDAR data. While the correlation module compute convolutional cross-correlation between the feature responses of adjacent key frames to 061 estimate displacement of the same objects. With local displacement and object orientation, 062 object location in intermediate frames can be calculated by interpolation. Moreover, detections can be linked between frames with the help of local displacement and multiple object tracking can be performed through *tracking by detection* [14].

In summary, our contributions are threefold: (i) we set up a two stream architecture based 066 on AVOD for simultaneous 3D object detection and tracking; (ii) we introduce correlation 067 features to capture object co-occurrences across time and perform frame-level detections 068 in a high speed through interpolation; (iii) we utilize tracking result to improve detection 069 performance and preliminary explore the algorithm for key frame selection.

074

076

077

087 088

091

## 2 Related Work

# 2.1 Video object detection

Video object detection in image has received increased attention since ImageNet VID datasets introduced. Most approaches for video object detection utilize optical flows, which present temporal information in videos. Some representative work such as FGFA [25], it leverages temporal coherence on feature level, and improves the pre-frame features by aggregation of nearby features along the motion paths with the help of optical flows. Later a more efficient approach based on [26] has been presented in [27], it introduces three new techniques: sparsely recursive feature aggregation, spatially-adaptive partial feature updating and temporally-adaptive key frame scheduling, which make this unified approach faster, more accurate and more flexible.

There are also some approach try to learn temporal information between consecutive frames. D&T [1] set up a ConvNet architecture for simultaneous detection and tracking in video. In order to capture cross-occurrences across time, it aid a correlation operation in networks. Our Bi-AVOD architecture mainly inspired by this work.

### 2.2 3D object detection

Currently, most approaches in 3D object detection can be divided into three types: image based detectors, LiDAR based detectors and fusion based detectors. Image based approaches such as Mono3D[1], 3DOP[1] use camera data only, since image has limited depth information, specific hand-crafted geometric features are required. LiDAR based methods are usually done in two fronts, one is utilizing a voxel grid representation to encoder point cloud and applying 3D CNN for features extracted, these approaches including 3D FCN [12], Vote3Deep [1] and VoxelNet [12] et al., these approaches suffer from the sparsity of point cloud and enormous computation cost in 3D convolution; others LIDAR based methods try to project point cloud to bird eye view (BEV) and apply 2D CNN for object detection, such as PIXOR[12], FaF[13] and Comple-YOLO [12] et al. These methods take advantage of the fact that objects in autonomous driving almost on the same plane thus loss of height information has little affect to the result, while the depth and Geometric information can be retained and computational complexity reduced significantly, making real-time detection possible. However, due to the sparsity of point cloud, the feature information after projecting is insufficient for accurate object detection especially for the small target.

There are also many multi-modal fusion methods that combine images and LiDAR data to improve detection accuracy. F-PointNet [27] first extracts the 3D bounding frustum of an object by extruding 2D bounding boxes from image detectors, then consecutively perform 3D object instance segmentation and amodal extent regression to estimate the amodal 3D bounding box. MV3D[2] extends the image based RPN of Faster R-CNN[27] to 3D and proposes a 3D RPN targeted at autonomous driving scenarios. MV3D uses every pixel in BEV feature map to multiple 3D anchors and then feeds the anchor to RPN to generate 3D proposals that are used to create view-specific feature crops from BEV feature maps and images. A deep fusion scheme is used to combine information from these feature crops to produce final detection output. However, MV3D does not work well for small targets due to the insufficient data for feature extracting caused by downsampling in convolutional feature extractors. AVOD[17] architecture is similar to MV3D in 3D RPN and feature fusion, however, its feature extract provides full resolution feature maps thus show greatly help in localization accuracy for small targets during the second stage of the detection framework. Our proposed architecture mostly based on AVOD mention above.

# 2.3 3D object tracking

More and more work has been done in 3D object tracking based on tracking by detection due to the rapidly development in 3D object detection. These approaches usually trend to apply 3D object detection and tracking simultaneously. FaF [13] jointly reasons about 3D detection, tracking and motion forecasting taking a 4D tensor created from multiple consecutive temporal frames. The most similar approach to our work is [13], however, their 3D detector is based on MV3D while ours is AVOD, and their detections association is done by solving a linear program after passing to a matching net and scoring net, while ours use a extending IOU based algorithm [2] by leveraging corresponding displacements over time.

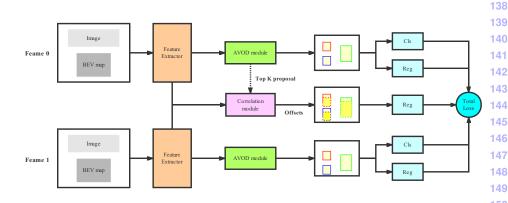


Figure 1: Bi-AVOD architecture

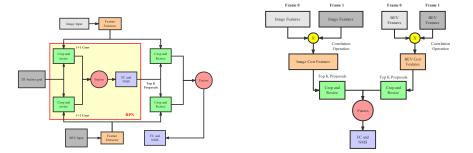


Figure 2: AVOD architecture

Figure 3: Correlation module

151

156

158

166

169

171

172 173

174

177

178

180

#### 3 **Methodology**

In this section we first give an overview of the Bi-AVOD approach (Sec. 3.1) that generates 168 detections and tracklets given two adjacent key frames as input. We then introduce the correlation module (Sec. 3.2) that aiding the network in the tracking process. Sec. 3.3 shows the multi-task objective function and Sec. 3.4 shows how we implement 3D streaming detection and tracking using prediction results of Bi-AVOD.

#### 3.1 **Bi-AVOD** model structure

We aim at performing 3D object detection in point cloud flow, and later apply multi-object tracking using tracking by detection paradigm. We build our Bi-AVOD mainly based on AVOD[5], Figure. illustrates our Bi-AVOD architecture. By doubling its input field, we can feed two adjacent key frames simultaneously and obtained corresponding object detection results. Meanwhile, the local displacements can be estimated by computing convolutional cross-correlation between the feature responses of adjacent frames using correlation module. With the predicted bounding boxes of two adjacent key frames and their local displacements, we can implement interpolation algorithms to generate interval bounding boxes for 3D flow detection and develop data association algorithms for 3D object tracking. Note that two *Feature Extractor* and two *AVOD module* in Figure. 1 are parameter sharing, thus the increased computational cost comes only from the correlation module. We will simply review the AVOD architecture in this section, and the correlation module will leave to the next section.

AVOD architecture is proposed for 3D object detection in autonomous driving by aggregating front view image and bird eye view (BEV) feature maps (generated by LiDAR data). It is a two stage method and Figure.2 illustrates its architecture. Firstly, it uses feature pyramids based extractors to generate full resolution feature maps from both BEV map and RGB image. Both feature maps are then feed to fusion RPN to generate  $Top\ K$  non-oriented region proposals after applied a  $1\times 1$  convolutional layer for dimensionality reduction, note that the default anchors in image and BEV map are generated through 3D anchor grid projection, and ROI Pooling is implemented by  $Crop\ and\ Resize$  operations. Finally, the  $Top\ K$  proposals are passed to the detection network for dimension refinement, orientation estimation and category classification. We refer the reader to [ $\Box$ ] for an further explanation of the architecture.

#### 3.2 Correlation module

The correlation module is illustrated in Figure.3. It first performs correlation on two image features and two BEV map features respectively. Then similar to the above mentioned RPN, it extracts feature crops via multiview *Crop and Resize* operations guided by *Top K* non-oriented region proposals. The feature crops then feed to fusion module and multiview aggregated features are generated. While the fusion module is identical to the one mentioned in AVOD, which first introduced in MV3D[ $\square$ ], and includes three different fusion schemes, *early fusion, late fusion* and *deep fusion*. Finally, the aggregated features are passed to a fully connected layer for regression after NMS.

Similar to FlowNet [ $\square$ ], we restrict correlation to a local neighborhood instead of all possible circular shifts in a feature map, which avoids large output dimensionality and too large displacements. The correlation operation performs point-wise feature comparison of two feature maps  $f_t$ ,  $f_{t+\tau}$  by

$$C^{t,t+\tau}(i,j,p,q) = \left\langle f_t(i,j), f_{t+\tau}(i+p,j+q) \right\rangle \tag{1}$$

where  $p,q\in[-d,d]$  are offsets to compare features in a local square window defined by the maximum displacement d, and i,j are the location of window center in feature map. The output is a correlation feature map of size  $\mathcal{C}\in\mathbb{R}^{h\times w\times (2d+1)\times (2d+1)}$  where h,w are the height, width of the feature map.

After above operation we get two correlation feature maps, one for point cloud  $\mathcal{C}_{pc}^{t,t+\tau}$ , and another for RGB image  $\mathcal{C}_{img}^{t,t+\tau}$ . Later ROI pooling and *early fusion* are performed just as in detection part. Aggregate feature  $\mathcal{C}_{fusion}^{t,t+\tau} = \frac{1}{2}(\mathcal{C}_{pc}^{t,t+\tau} + \mathcal{C}_{img}^{t,t+\tau})$  is then flatten and fed to a fully connected layer to predict the transformation  $\Delta^{t,t+\tau} = (\Delta_{x,y,z}^{t,t+\tau}, \Delta_{w,h,l}^{t,t+\tau}, \Delta_r^{t,t+\tau}) \in \mathbb{R}^7$  of the RoIs from t to  $t+\tau$ .

## 3.3 Multitask detection and correlation objective

We extend the multi-task loss of object detection, consisting of a classification loss  $L_{cls}$  and a regression loss  $L_{reg}$ , with an additional term  $L_{corr}$  that scores the displacement regression be-

tween objects across two frames. Considering a batch of N RoIs after category balanced sam- 230 pling, the network predicts softmax probabilities  $\{p_i\}_{i=1}^N$ , bounding box regression offsets 231  $\{b_i\}_{i=1}^N$ , and cross-frame displacement regression  $\{\Delta_i^{t+\tau}\}_{i=1}^{N_{corr}}$ , the overall objective function 232 is shown as:

$$L(\{p_i\},\{b_i\},\{\Delta_i\}) = \frac{1}{N} \sum_{i=1}^{N} L_{cls}(p_{i,c^*}) + \frac{\alpha}{N_{fg}} \sum_{i=1}^{N} [c_i^* > 0] L_{reg}(b_i,b_i^*)$$

$$+ \frac{\beta}{N_{corr}} \sum_{i=1}^{N_{corr}} L_{corr}(\Delta_i^{t+\tau},\Delta_i^{*,t+\tau})$$

$$(2) \frac{236}{238}$$

$$(3) \frac{236}{238}$$

237

241

243

245

246

247

249 250 251

252

254

256

257

260

264

267

269

274

275

where  $c_i^*$  is the ground truth class label of an RoI and  $p_{i,c_i^*}$  is corresponding predicted softmax score.  $b_i^*$  is the ground truth bounding box regression target, and  $\Delta_i^{*,t+\tau}$  is the displacement regression target. The indicator function  $[c_i^* > 0]$  is 1 for foreground RoIs and 0 for background RoIs.  $L_{cls}$ ) is cross-entropy loss, and  $L_{reg}$ ,  $L_{corr}$  are smooth L1 loss [ $\square$ ].  $\alpha$  and  $\beta$  are weight for  $L_{reg}$  and  $L_{corr}$ . Note that we only consider the  $N_{fg}$  foreground RoIs loss for  $L_{reg}$ , and  $L_{corr}$  is active only for foreground RoIs which have a track correspondence across two key frames. Additionally, For a single target  $D_t = (D_t^{x,y,z}, D_t^{w,h,l}, D_t^r) \in \mathbb{R}^7$  in frame t and its track correspondence  $D_{t+\tau}$  in frame  $t+\tau$ , the displacement regression values for the target  $\Delta_i^{*,t+\tau}$  are encoded just like the 3D bounding box in [ $\square$ ].

## 3D streaming detection and tracking

Give a streaming point cloud of N frames  $\{I_f\}$  for  $f \in \{1,...N\}$ , the streaming object detection task needs to predict a set of detections  $D_f$  for each frame  $I_f$ . Each detection set  $D_f$ consists of object detections  $\{D_f^i\}$  while  $i \in \{1,...N_f\}$   $(N_f$  is the number of detections in frame f). Note that  $D_f$  can also be an empty set when no object is detected in a frame. In 3D object detection, Each detection  $D_f^l$  is parametrized as  $D_f^l = (x_f^l, y_f^l, z_f^l, w_f^l, h_f^l, l_f^l, \theta_f^l, s_f^l)$ , where  $(x_f^i, y_f^i, z_f^i)$  corresponds to the center (bottom center in KITTI[ $\square$ ] Datasets) of the detection box in point cloud,  $(w_f^i, h_f^i, l_f^i)$  corresponds to width, height, length of the detection box,  $\theta_f^i$  is the rotation angle in yaw axis and  $s_f^i$  is the detectors confidence in the bounding 261 box.

Because of the redundant features in streaming, we can only compute the object detections in key frames, while the detections in intermediate frames can be calculated using the detections in adjacent two key frames. Support we have two predicted object detections set  $(D_f, D_{f+\tau})$  in two consecutive key frames  $(I_f, I_{f+\tau})$  where  $\tau$  is temporal stride, the detection results  $\{D_{f+t}^i\}$   $(t \in \{1, \tau - 1\})$  in intermediate frame  $D_{f+t}$  can be obtained by:

$$D_{f+t}^i = \mathcal{F}(W_f D_f^i, W_{f+\tau} D_{f+\tau}^i) \tag{3}$$

where W is corresponding weight and  $\mathcal{F}$  is generation function to produce  $D_{f+t}^i$ , in this 270paper we use quadratic interpolation function. Note that we can utilize (1) to generate  $D_{f+t}^i$ only when the target exists in  $D_f$  and  $D_{f+\tau}$  simultaneously. if the target is emerged or end in intermediate frames, this method would be failed. Though we can develop the key frames selection function carefully to handle this situation, it is beyond the scope of this article. In this paper we focus on the targets that always exist between two key frames.

With object detections in each frame, multi-object tracking can be implemented by tracking-by-detection paradigm. For each bounding box in each frame, MOT try to associate it to a unique target trajectory  $T_k = \{D_{f_1}^k, D_{f_2}^k, ..., D_{f_{N_k}}^k\}$ , where k is trajectory id and  $N_k$  is the length of  $T_k$ ,  $\{f_1, f_2, ..., f_{N_k}\}$  are corresponding frames id.

## 4 Experimental evaluation

#### 4.1 Datasets

276

277

282283

287

295

297

301

311

313

314

315 316

317

318

319

321

We use the KITTI object tracking Benchmark [LL] for evaluation. It consists of 21 training sequences and 29 test sequences with vehicles annotated in 3D. Each sequence includes hundreds of point clouds frames captured by Velodyne HDL-64E rotating 3D laser scanner and corresponding RGB images. We split 21 training sequences into two parts based on the parity of the sequence number, odd number for training and even number for evaluation. For testing ,we train our model in all 21 training sequences.

## 4.2 Experiments

## 5 Conclusions

#### References

- [1] Aseem Behl, Despoina Paschalidou, Simon Donné, and Andreas Geiger. Pointflownet: Learning representations for 3d scene flow estimation from point clouds. *arXiv preprint arXiv:1806.02170*, 2018.
- [2] Erik Bochinski, Tobias Senst, and Thomas Sikora. Extending iou based multi-object tracking by visual information. *AVSS. IEEE*, 2018.
- [3] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun. Monocular 3d object detection for autonomous driving. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 2147–2156, June 2016. doi: 10.1109/CVPR.2016. 236.
- [4] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1907–1915, 2017.
- [5] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals using stereo imagery for accurate object class detection. *IEEE transactions on pattern analysis and machine intelligence*, 40(5):1259–1272, 2018.
- [6] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016.
- [7] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet:

[8] Martin Engelcke, Dushyant Rao, Dominic Zeng Wang, Chi Hay Tong, and Ingmar Pos-

national conference on computer vision, pages 2758–2766, 2015.

Learning optical flow with convolutional networks. In Proceedings of the IEEE inter- 322

323324

	ner. Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks. In 2017 IEEE International Conference on Robotics and Automation (ICRA), pages 1355–1361. IEEE, 2017.	<ul><li>326</li><li>327</li><li>328</li></ul>
[9]	Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Detect to track and track to detect. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , pages 3038–3046, 2017.	<ul><li>329</li><li>330</li><li>331</li><li>332</li></ul>
[10]	Davi Frossard and Raquel Urtasun. End-to-end learning of multi-sensor 3d tracking by detection. In 2018 IEEE International Conference on Robotics and Automation (ICRA), pages 635–642. IEEE, 2018.	333 334 335
[11]	Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. <i>The International Journal of Robotics Research</i> , 32(11): 1231–1237, 2013.	<ul><li>336</li><li>337</li><li>338</li><li>339</li></ul>
[12]	Ross Girshick. Fast r-cnn. In <i>Proceedings of the IEEE international conference on computer vision</i> , pages 1440–1448, 2015.	340 341
[13]	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 770–778, 2016.	<ul><li>342</li><li>343</li><li>344</li><li>345</li></ul>
[14]	Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In <i>Advances in neural information processing systems</i> , pages 1097–1105, 2012.	<ul><li>346</li><li>347</li><li>348</li><li>349</li></ul>
[15]	Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 1–8. IEEE, 2018.	
[16]	Philip Lenz, Andreas Geiger, and Raquel Urtasun. Followme: Efficient online mincost flow tracking with bounded memory and computation. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , pages 4364–4372, 2015.	354 355 356 357
[17]	Bo Li. 3d fully convolutional network for vehicle detection in point cloud. In 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 1513–1518. IEEE, 2017.	359 360
[18]	Xingyu Liu, Charles R Qi, and Leonidas J Guibas. Learning scene flow in 3d point clouds. <i>arXiv preprint arXiv:1806.01411</i> , 2018.	<ul><li>361</li><li>362</li><li>363</li></ul>
[19]	Wenjie Luo, Bin Yang, and Raquel Urtasun. Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 3569–3577, 2018.	364 365 366 367

371

374

379

391

392

- 368 [20] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets 369 for 3d object detection from rgb-d data. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 918–927, 2018. 370 [21] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real
  - time object detection with region proposal networks. In Advances in neural information processing systems, pages 91–99, 2015.
  - [22] Martin Simon, Stefan Milz, Karl Amende, and Horst-Michael Gross. Complex-volo: An euler-region-proposal for real-time 3d object detection on point clouds. In European Conference on Computer Vision, pages 197–209. Springer, 2018.
  - [23] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for largescale image recognition. arXiv preprint arXiv:1409.1556, 2014.
  - [24] Bin Yang, Wenjie Luo, and Raquel Urtasun. Pixor: Real-time 3d object detection from point clouds. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 7652–7660, 2018.
  - [25] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 4490–4499, 2018.
  - [26] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. Flow-guided feature aggregation for video object detection. In Proceedings of the IEEE International Conference on Computer Vision, pages 408–417, 2017.
  - [27] Xizhou Zhu, Jifeng Dai, Lu Yuan, and Yichen Wei. Towards high performance video object detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 7210–7218, 2018.