

---

# YOLO4D: A Spatio-temporal Approach for Real-time Multi-object Detection and Classification from LiDAR Point Clouds

---

**Ahmad El Sallab\***

Valeo Egypt AI Research  
ahmad.el-sallab@valeo.com

**Ibrahim Sobh\***

Valeo Egypt AI Research  
ibrahim.sobh@valeo.com

**Mahmoud Zidan\***

Valeo Egypt AI Research  
mahmoud.ismail-zidan.ext@valeo.com

**Mohamed Zahran\***

Valeo Egypt AI Research  
mohamed.zahran@valeo.com

**Sherif Abdelkarim\***

Valeo Egypt AI Research  
sherif.abdelkarim@valeo.com

## Abstract

In this paper, YOLO4D is presented for Spatio-temporal Real-time 3D Multi-object detection and classification from LiDAR point clouds. Automated Driving dynamic scenarios are rich in temporal information. Most of the current 3D Object Detection approaches are focused on processing the spatial sensory features, either in 2D or 3D spaces, while the temporal factor is not fully exploited yet, especially from 3D LiDAR point clouds. In YOLO4D approach, the 3D LiDAR point clouds are aggregated over time as a 4D tensor; 3D space dimensions in addition to the time dimension, which is fed to a one-shot fully convolutional detector, based on YOLO v2 architecture. The outputs are the oriented 3D Object Bounding Box information, together with the object class. Two different techniques are evaluated to incorporate the temporal dimension; recurrence and frame stacking. The experiments conducted on KITTI dataset, show the advantages of incorporating the temporal dimension.

## 1 Introduction

Environmental modeling and perception is a critical component to automated driving pipeline. In order to have efficient planning in complex scenarios, 3D object detection is essential, to get information about the objects extent and range in the 3D space. Deep Learning is becoming the trend for real-time object detection and classification [13][15] [1] [14][9]. LiDAR sensors are used to perceive the 3D nature of objects, where the sensor provides a 3D point cloud (PCL) representing the range of reflected laser beams of the surrounding objects. The task of 3D object detection and classification from such a point cloud is challenging. On one hand, the point cloud is sparse, since not all beams are reflected, and noisy due to imperfect reflections and echoes. On the other hand, the 3D point cloud does not have the color and texture features that characterize the object classes as in the case of 2D camera perspective images. Such complexity, in addition to the dynamic nature of the environment, motivates our work to incorporate the temporal factor in addition to the spatial features of the input 3D LiDAR point clouds.

---

\*The authors are ordered alphabetically

In this paper, we present YOLO4D; a Spatio-temporal extension of the work done in YOLO3D[1] for real-time multi-object detection and classification from 3D LiDAR point clouds, where YOLO3D is extended with Convolutional LSTM [18] for temporal features aggregation. The 3D LiDAR point clouds are aggregated over time as a 4D tensor; 3D space dimensions in addition to the time dimension. The output is the oriented 3D object bounding box (OBB) center coordinates, in addition to its length (L), width (W), height (H) and orientation (yaw), together with the objects classes and confidence scores. The evaluation of the Spatio-temporal approach is performed on KITTI dataset [4]. We compare two approaches to incorporate the temporal factor; YOLO4D and frame stacking. The experimental results show the advantage of adding the temporal aggregation in addition to spatial features.

The rest of the paper is organized as follows; first, we discuss the related work, followed by the Spatio-temporal approach, with the proposed network architecture. Finally, we present the experimental results and evaluation of different techniques on KITTI dataset.

## 2 Related Work

**3D Object Detection:** Recent works in 3D object detection depend on 3D sensors like LiDAR to take advantage of accurate depth information. Different data representations are introduced, like projecting point cloud in 2D view (Bird Eye View, Front View) such as PIXOR [19] , [8] and MV3D [3]. PIXOR assumes that all objects lie on the same ground. In [2] the point cloud is converted to a front view depth map. Others like [7] and [20], convert the point cloud to voxels. These LiDAR 3D point clouds object detection methods do not take the advantage of temporal information to produce more accurate 3D bounding boxes.

**Spatial-temporal Object Detection:** While some papers addressed the use of recurrence as Spatial-temporal feature extractor alongside with the object detection [17] [12], they focused on improving visual tracking in videos. For example, in ROLO [12] the same architecture of YOLO v1 [13] is used with an LSTM [5] layer added at the end. The network consumes as input raw videos and returns 2D tracked bounding box. Recently, Fast and Furious [10] incorporates the time with 3D voxels using 3D convolutions and adopts a multi-task learning setup.

In our work we extend YOLO3D [1] to detect the 3D bounding boxes from a single channeled bird eye view. Additionally, a Convolutional LSTM layer is injected to encode the temporal information.

## 3 Spatio-temporal 3D object detection approach

In this section, the approach for Spatio-temporal 3D object detection is described. The main intuition behind our work is to leverage not only the spatial but also the temporal information in input sequences for more accurate object detection. For encoding temporal sequences, we adopted two different approaches, YOLO4D and frame stacking. These approaches encode the temporal information in different ways. Frame stacking is a simple method that depends on the input layer, while, YOLO4D employs Convolutional LSTM.

### 3.1 Spatial 3D object detection

Following the work done in YOLO3D [1], the point cloud is projected into a bird's eye view (BEV) grid map. The orientation angle of the bounding boxes is normalized and used as a single regressed value. For 3D bounding box regressions, two regression terms are added to the original YOLO architecture, the  $z$  coordinate of the center, and the height  $h$  of the box. The average 3D box dimensions for each object class is calculated on the training dataset and used as anchors. The loss for the 3D oriented boxes is an extension to the original YOLO loss for 2D boxes. The total loss is formulated as the weighted summation of the mean squared error over the 3D coordinates, dimensions, the mean squared error over the angle, the confidence score, and the cross-entropy loss over the object classes.

### 3.2 Temporal aggregation

The dataset  $\mathcal{T}$  is composed of number of  $k$  scenarios  $\{S^0, S^1, \dots, S^{k-1}\}$ , a scenario  $S^i$  consists of a sequence of  $n$  frames  $\{I_0^i, I_1^i, \dots, I_{n-1}^i\}$ , with the corresponding list of 3D bounding boxes targets for each frame,  $\{D_0^i, D_1^i, \dots, D_{n-1}^i\}$ . Each scenario is divided into a number of short clips of  $m$  frames. The frames of the clip  $j$  from scenario  $i$  is defined as  $Q_j^i = \{I_{t_j}^i, I_{t_j+1}^i, \dots, I_{t_j+m-1}^i\}$ .

#### 3.2.1 Frame stacking

In frame stacking, frames of each clip are stacked in order together and presented as a single input to the object detection network for training. Accordingly, the frame stacking approach encodes the temporal information indirectly by reshaping the input by increasing its depth to represent the changes over time. During the training process, it is up to the network to learn the temporal information from the input stacked frames without encoding hidden state through recurrent layers. The loss is similar to the YOLO3D architecture for a single frame input, where the ground truth bounding boxes of the last frame are used. Accordingly, the predicted 3D bounding boxes depend only on a single stacked input, as shown in Figure 1 (left).

#### 3.2.2 Convolutional LSTM

In YOLO4D architecture, a Convolutional LSTM layer is injected directly into YOLO3D single frame architecture, see section 4.3 for more details. Convolutional LSTM allows the network to learn both spatial and temporal information. The network is trained on the same sequences used for the frame stacking approach, leading to a model that is capable of detecting objects in temporal streams of input frames. **The prediction model can be considered as a function  $F$** , parameterized by  $\theta$ , that maps an input frame  $I$  and the previous state  $s_{t-1}$  to a list of 3D bounding boxes  $D$ , as shown in Eq.(1).

$$\mathcal{F}_\theta(I_t, s_{t-1}) = (D_t, s_t) \quad (1)$$

Where the state  $s_t$  is used as input for the next time step predictions. The loss in this case is the same loss of YOLO3D however, the optimization is back-propagated through time via the injected Convolutional LSTM layer to maintain the temporal information.

Figure 1, illustrates the frame stacking and Convolutional LSTM object detection architectures.

Because of the recurrency, the network generated values depend not only on the current frame at time  $t$ , but also on the previous  $m$  frames. On the other hand, the ground truth values correspond to the current frame. The ground truth values are defined as  $v^{(t)}$ , while the network generated values at time

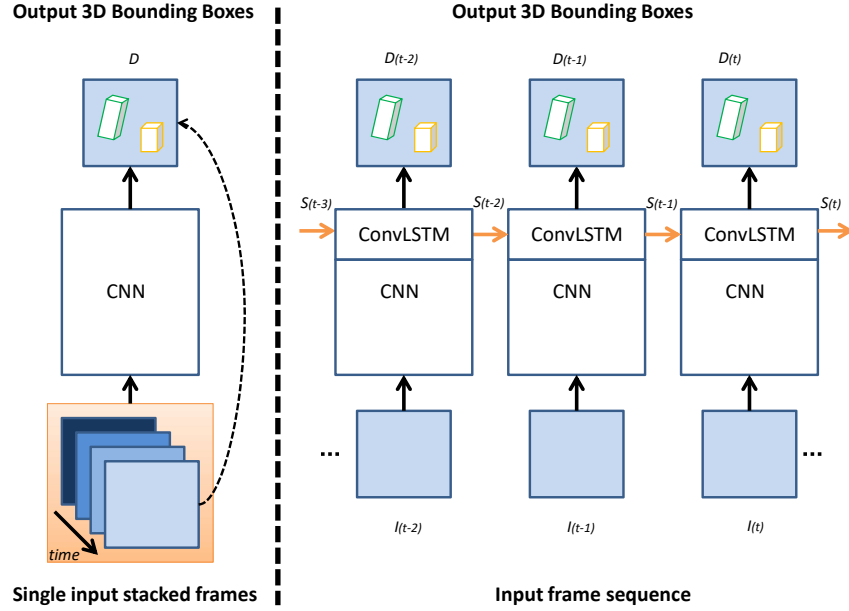


Figure 1: **Left:** Frame stacking architecture; **Right:** Convolutional LSTM architecture.

$t$  are defined as  $\hat{v}^{(t_m)}$ , where  $t_m$  indicates the value generated based on input  $m$  frames, from time  $t$  back to time  $t - m - 1$ . The total loss is calculated as shown in Eq.(2).

$$\begin{aligned}
L_{\theta} = & \lambda_{coord} \sum_{i=0}^{s^2} \sum_{j=0}^B L_{ij}^{obj} \left( (x_i^{(t)} - \hat{x}_i^{(t_m)})^2 + (y_i^{(t)} - \hat{y}_i^{(t_m)})^2 + (z_i^{(t)} - \hat{z}_i^{(t_m)})^2 \right) \\
& + \lambda_{coord} \sum_{i=0}^{s^2} \sum_{j=0}^B L_{ij}^{obj} \left( (\sqrt{w_i^{(t)}} - \sqrt{\hat{w}_i^{(t_m)}})^2 + (\sqrt{l_i^{(t)}} - \sqrt{\hat{l}_i^{(t_m)}})^2 + (\sqrt{h_i^{(t)}} - \sqrt{\hat{h}_i^{(t_m)}})^2 \right) \\
& + \lambda_{yaw} \sum_{i=0}^{s^2} \sum_{j=0}^B L_{ij}^{obj} (\phi_i^{(t)} - \hat{\phi}_i^{(t_m)})^2 \\
& + \lambda_{obj} \sum_{i=0}^{s^2} \sum_{j=0}^B L_{ij}^{obj} (C_i^{(t)} - \hat{C}_i^{(t_m)})^2 \\
& + \lambda_{noobj} \sum_{i=0}^{s^2} \sum_{j=0}^B L_{ij}^{noobj} (C_i^{(t)} - \hat{C}_i^{(t_m)})^2 \\
& + \lambda_{classes} \sum_{i=0}^{s^2} \sum_{j=0}^B L_{ij}^{obj} \left( - \sum_{c \in classes} p_i^{(t)}(c) \log(\hat{p}_i^{(t_m)}(c)) \right)
\end{aligned} \tag{2}$$

Where  $\lambda_{coord}$  is the weight assigned to the loss over the coordinates.  $\lambda_{obj}, \lambda_{noobj}$  are the weights assigned to the loss over predicting the confidences for objects and no objects respectively.  $\lambda_{yaw}$  is the weight assigned to the loss over the orientation angle.  $\lambda_{classes}$  is the weight assigned to the loss

over the class probabilities.  $L_{ij}^{obj}$  is a variable that takes the value of 1 if there is a ground truth box in the  $i$ th cell and the  $j$ th anchor is the associated anchor, otherwise 0.  $L_{ij}^{noobj}$  is the complement of the previous variable, takes the value of 1 if there is no object, and 0 otherwise.  $x_i, y_i, z_i$  : are the ground truth coordinates.  $\hat{x}_i, \hat{y}_i, \hat{z}_i$  are the predicted coordinates.  $\phi_i, \hat{\phi}_i$  are the ground truth and predicted orientation angle respectively.  $C_i, \hat{C}_i$  are the ground truth and predicted confidence respectively.  $w_i, l_i, h_i$  are the ground truth width, length, and height of the box respectively.  $\hat{w}_i, \hat{l}_i, \hat{h}_i$  are the predicted width, length, and height of the box respectively.  $p_i(c), \hat{p}_i(c)$  are the ground truth and predicted class probabilities respectively.  $B$  is the number of boxes, and  $s$  is the length of one of the sides of the square output grid, thus  $s^2$  is the number of grids in the output.

## 4 Experimental setup

### 4.1 Dataset

All the experiments are conducted on the publicly available KITTI raw dataset [4], which consists of sequenced frames, unlike the KITTI benchmark dataset [4]. The dataset consists of 36 different annotated point cloud scenarios of variable lengths and a total of 12919 frames. These scenarios are divided into clips as described in section 3.2, with  $m = 4$ . Moreover, these scenarios have diverse driving environments such as highway roads, traffic, city and residential areas. They are also rich with dynamic and static objects.

We study the effect of Spatio-temporal object detection on 5 classes of objects. Our objects of interest are Pedestrians, Cyclists, Cars, Vans, and Trucks. 80% of the dataset is used for training, and the remaining 20% for evaluation. We choose not to ignore Vans or Trucks or devise a separate model for each class, as popular approaches in LiDAR-based object detection do [6][3]. In contrast, we are benefiting from YOLO as a single shot detector recognizing all KITTI classes in the same time, which makes our approach more practical for automated driving.

### 4.2 Bird Eye View (BEV)

The same point cloud pre-processing procedure in YOLO3D [1] is followed to generate the BEV input. In particular, we project the PCL into a single channel height grid map of size  $608 \times 608$ , at a cell resolution of 0.1m. The height grid map encodes the height of the highest PCL point associated with that cell. The major difference from YOLO3D’s BEV is using a single channeled input instead of two, that is, we discarded the density channel, because we found that training with the height channel only did not significantly hurt the performance, while reducing the memory footprint, and allows more efficient training.

### 4.3 Architectures

Two fully convolutional object detection architectures are adopted for the experiments: Mixed-YOLO, which is YOLO-v2[15] in half precision, following [11] for a mixed precision training and Tiny-YOLO [14] in full precision [13]. We extended them for oriented 3D bounding boxes detection as mentioned in section 3.1, hence we call them Mixed-YOLO3D and Tiny-YOLO3D respectively. The motivation behind using mixed precision training and having the model weights in half precision is because that YOLO4D full precision model requires very high memory footprint to train effectively with a reasonable batch size. Tiny-YOLO [14] is adopted to experiment the Spatio-temporal effect on shallower models.

**Frame stacking:** as described in section 3.2.1, each clip’s  $m$  frames are stacked along the channel dimension and the input layer of Mixed-YOLO3D and Tiny-YOLO3D is modified accordingly.

**Convolutional LSTM:** as described in section 3.2.2, Mixed-YOLO3D and Tiny-YOLO3D networks are extended to account for the temporal dimension, by injecting a Convolutional LSTM layer just before the final output layer, revealing Mixed-YOLO4D and Tiny-YOLO4D respectively. This recurrent layer takes the  $19 \times 19 \times 1024$  feature map, produced by the last convolutional hidden layer in both models as an input per time step, and outputs 512 channels using a  $3 \times 3$  kernel size, which encodes the Spatio-temporal features that are fed to the final output layer for object detection and classification.

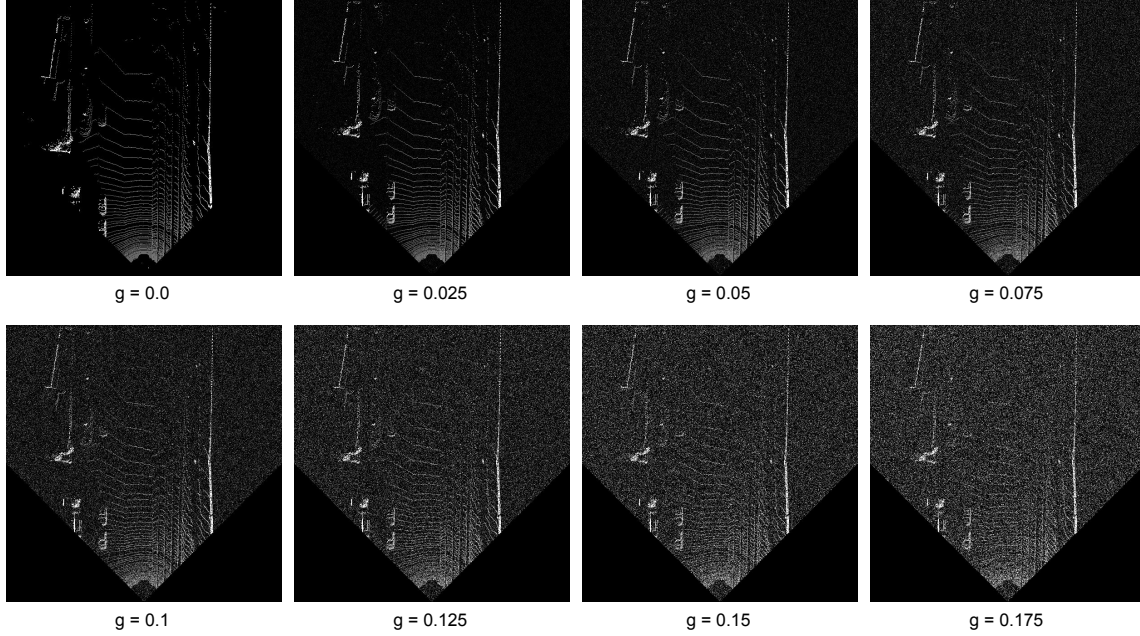


Figure 2: BEV inputs at different Gaussian noise scales  $g$

#### 4.4 Training

For all Spatio-temporal based models, a clip length  $m$  of 4 is used. All models are trained until convergence, with a fixed batch size of 4, and a weight decay of  $5e - 5$ . For optimization, stochastic gradient descent (SGD) is used with a momentum of 0.99, and a learning rate of  $1e - 4$ . Regarding the half precision models: Mixed-YOLO3D, Mixed-YOLO4D and Frame Stacking + Mixed-YOLO3D, we followed [11] for a mixed precision training, forward pass in half precision, loss computation and weights update in full precision. No loss scaling is used for Mixed-YOLO3D and a loss scaling of 8 is used for the other 2 models.

#### 4.5 Robustness

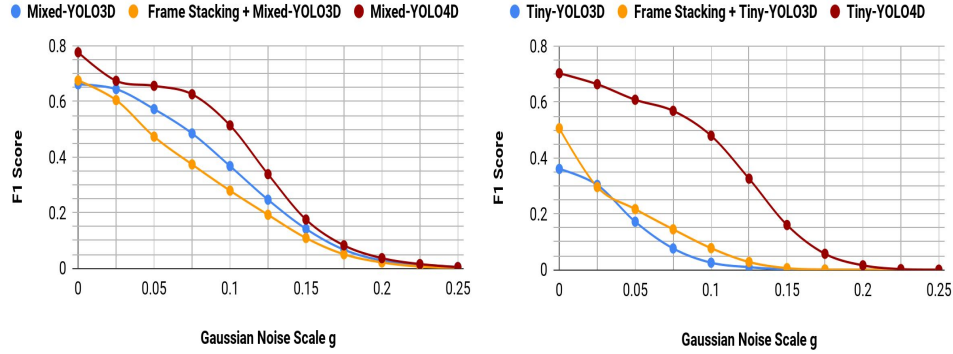
In real-world applications, it is important to evaluate the performance in case of failure in sensory readings, for example, [16] adopted frame dropping. In this work, noisy inputs are conducted by adding Gaussian noise to the BEV input frames at different scales  $g$ , then the detection performance is examined. Examples of inputs at different Gaussian noise scales are shown in Figure 2.

### 5 Results

The performance of frame stacking and YOLO4D models are compared to assess the effect of temporal information. As a baseline, Mixed-YOLO3D and Tiny-YOLO3D, are also compared. The F1 scores on the validation set are shown in Table 1. Based on the experimental results, the deeper Mixed-YOLO models achieve better mean F1 scores compared to the shallower Tiny-YOLO counterparts. As expected, the YOLO4D models outperform the frame stacking models. Frame stacking encodes the temporal information only through the reshaping of inputs, while YOLO4D encodes the temporal information in a more natural way through the recurrent convolutional LSTM layer allowing a better propagation of temporal representations through time. YOLO4D models outperform all other methods on all classes, achieving 11.5% improvement on Mixed-YOLO3D, and 34.26% improvement on Tiny-YOLO3D. Frame stacking provides a 2.36% improvement on the Mixed-YOLO baseline model, and a 14.56% improvement on the shallower Tiny-YOLO baseline model. Mixed-YOLO4D achieves the best mean F1 score, taking 20ms on TITAN XP GPU.

Table 1: Performance Comparisons: F1 Score

Model	Mean F1	Car	Pedestrian	Cyclist	Truck	Van
<b>Mixed-YOLO:</b>						
Mixed-YOLO3D	66.23%	69.82%	5.59%	14.61%	68.69%	64.32%
Frame Stacking+Mixed-YOLO3D	68.59%	71.77%	6.42%	16.93%	71.62%	65.45%
Mixed-YOLO4D	<b>77.73%</b>	<b>82.16%</b>	<b>10.44%</b>	<b>27.19%</b>	<b>81.77%</b>	<b>80.88%</b>
<b>Tiny-YOLO:</b>						
Tiny-YOLO3D	36.10%	37.77%	1.53%	2.29%	39.17%	36.26%
Frame Stacking+Tiny-YOLO3D	50.66%	53.69%	2.56%	2.86%	60.32%	43.13%
Tiny-YOLO4D	<b>70.36%</b>	<b>74.24%</b>	<b>10.03%</b>	<b>19.49%</b>	<b>73.00%</b>	<b>69.53%</b>

Figure 3: Robustness Performance: **Left:** Mixed-YOLO; **Right:** Tiny-YOLO

As shown in Figure 3, YOLO4D models exhibit more robustness across different scales of noisy inputs. The performance of YOLO3D baseline models and frame stacking models drop significantly, while YOLO4D models maintain almost the same performance between noise scales of 0.025 and 0.075. The frame stacking and baseline models have comparable performances. As expected, Frame Stacking + Tiny-YOLO3D model shows slightly more robustness compared to the baseline Tiny-YOLO3D. On the other hand, the baseline Mixed-YOLO3D shows slightly better robustness than the Frame Stacking + Mixed-YOLO3D model.

## 6 Conclusion

In this work, YOLO4D is proposed for Spatio-temporal Real-time 3D Multi-object detection and classification from LiDAR point clouds, where the inputs are 4D tensors encoding the spatial 3D information and temporal information, and the outputs are the oriented 3D object bounding boxes information, together with the object class and confidence score. The experimental results show the effect of adding the temporal in addition to the spatial features for achieving better detection. Recurrence and frame stacking are evaluated to incorporate the temporal dimension on KITTI dataset. Both recurrence and frame stacking show better detection performance compared to single frame detection. However, as expected, recurrent YOLO4D achieves a better detection compared to frame stacking. Furthermore, robustness of the detection in the presence of noisy inputs is evaluated and it is clear that the YOLO4D models are more robust than frame stacking and single frame models.

## References

- [1] W. Ali, S. Abdelkarim, M. Zahran, M. Zidan, and A. E. Sallab. Yolo3d: End-to-end real-time 3d oriented object bounding box detection from lidar point cloud. *arXiv preprint arXiv:1808.02350*, 2018.
- [2] A. Asvadi, L. Garrote, C. Premevida, P. Peixoto, and U. J. Nunes. Depthcn: Vehicle detection using 3d-lidar and convnet. In *Intelligent Transportation Systems (ITSC), 2017 IEEE 20th International Conference on*, pages 1–6. IEEE, 2017.
- [3] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia. Multi-view 3d object detection network for autonomous driving. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6526–6534. IEEE, 2017.
- [4] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- [5] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [6] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. Waslander. Joint 3d proposal generation and object detection from view aggregation. *arXiv preprint arXiv:1712.02294*, 2017.
- [7] B. Li. 3d fully convolutional network for vehicle detection in point cloud. In *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*, pages 1513–1518. IEEE, 2017.
- [8] B. Li, T. Zhang, and T. Xia. Vehicle detection from 3d lidar using fully convolutional network. *arXiv preprint arXiv:1608.07916*, 2016.
- [9] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [10] W. Luo, B. Yang, and R. Urtasun. Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net.
- [11] P. Micikevicius, S. Narang, J. Alben, G. Diamos, E. Elsen, D. Garcia, B. Ginsburg, M. Houston, O. Kuchaev, G. Venkatesh, et al. Mixed precision training. *arXiv preprint arXiv:1710.03740*, 2017.
- [12] G. Ning, Z. Zhang, C. Huang, X. Ren, H. Wang, C. Cai, and Z. He. Spatially supervised recurrent convolutional neural networks for visual object tracking. In *Circuits and Systems (ISCAS), 2017 IEEE International Symposium on*, pages 1–4. IEEE, 2017.
- [13] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [14] J. Redmon and A. Farhadi. Tiny yolo.
- [15] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 6517–6525. IEEE, 2017.
- [16] I. Sobh and N. Darwish. Robust dual view deep agent. In *Proceedings of the 2nd International Sino-Egyptian Congress on Agriculture, Veterinary Sciences and Engineering*, 2017.
- [17] S. Tripathi, Z. C. Lipton, S. Belongie, and T. Nguyen. Context matters: Refining object detection in video with recurrent neural networks. *arXiv preprint arXiv:1607.04648*, 2016.
- [18] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, pages 802–810, 2015.
- [19] B. Yang, W. Luo, and R. Urtasun. Pixor: Real-time 3d object detection from point clouds.
- [20] Y. Zhou and O. Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. *arXiv preprint arXiv:1711.06396*, 2017.