

UNIVERSITÄT STUTTGART

MASTER THESIS

**3D Pose Estimation of Vehicles from Monocular
Videos using Deep Learning**

Qing Cheng

Course of Study: Information Technology

Examiner: Ph.D. Daniel Hennes

Supervisors: Dr. Hung Ngo, Jie Zhong

Commenced: January 1, 2018

Completed: July 3, 2018

Abstract

In this thesis, we present a novel approach to perform 3D pose estimation of vehicles from monocular images intended for autonomous driving scenarios. A robust deep neural network is applied to perform 3D dimension proximity estimation, 2D part localization, and 2D part visibility prediction simultaneously. In the inference phase, these learned features are fed to a pose estimation algorithm to recover the 3D location, 3D orientation, and 3D dimensions of the vehicles with the help of a set of 3D vehicle models. Our approach can perform these six tasks simultaneously in real time and handle highly occluded or truncated vehicles. The experiment results show that our approach achieves state-of-the-art performance on six tasks and outperforms most of monocular methods on the challenging KITTI benchmark.

Keywords: 3D pose estimation, 3D vehicle detection, deep learning, computer vision

Contents

1	Introduction	4
1.1	Motivation	4
1.2	Contributions	4
2	Related Work	6
2.1	Mono RGB Image Based Approaches	6
2.2	Stereo RGB Image Based Approaches	7
2.3	Point Cloud Based Approaches	7
2.4	Sensor Fusion Based Approaches	8
3	Background	10
3.1	A Short History of Autonomous Vehicles	10
3.2	Deep Learning Technology	11
3.3	Computer Vision	13
3.3.1	Image Formation	14
3.3.2	Intrinsic and Extrinsic Parameters	14
4	My Work - Technical Detail	17
4.1	Data and Labelling	18
4.1.1	KITTI Dataset	18
4.1.2	3D Vehicle Dataset	19
4.1.3	Labelling for KITTI Dataset	20
4.1.4	Notations	23
4.2	Deep CVT Network	23
4.2.1	Architecture	23
4.2.2	Training	28
4.3	Inference	34
4.3.1	Template Matching	34
4.3.2	2D-3D Matching	35
4.3.3	Visualization	35
5	Experiment and Evaluation	37
5.1	Experiment Setup	37
5.2	Evaluation Metrics	37
5.3	Design Choices	39
5.3.1	Model Selection Strategy	39
5.3.2	Difficulty Levels For Dataset Division	39
5.3.3	Cost Function	41

5.3.4	Learning Rate and Weight Decay	43
5.3.5	Feature Extractor	44
5.3.6	Selection of Characteristic Points	45
5.4	Experiment Results	47
5.4.1	3D Bounding Box Estimation	47
5.4.2	3D Localization	48
5.4.3	3D Orientation Estimation	49
5.4.4	3D dimension estimation, 2D part localization, and 2D part visibility	50
6	Discussion	51
6.1	Deficiency of 2D Coordinate Labelling	51
6.1.1	Deficiency of Our Labelling Approach	51
6.1.2	Deficiency of KITTI Dataset Labelling	53
6.1.3	Possible Improvements	53
6.2	Deficiency of Visibility Labelling	53
6.2.1	Deficiency of Our Labelling Approach	53
6.2.2	Possible Improvements	54
6.3	Deficiency of Template Matching	55
6.3.1	Possible Improvements	56
6.4	Deficiency of Input Image Format	56
6.4.1	Possible Improvements	57
6.5	Deficiency of the Whole Approach	57
7	Conclusion	59

1 Introduction

1.1 Motivation

The autonomous vehicles (AVs) is the evolutionary direction of automobiles due to its promising reliability and efficiency, as well as the underlying commercial profits. To gain a comprehensive perception of the driving environment is the very essential step for AVs. Object detection is one of the key challenges of perception. Thanks to the remarkable advancement of deep neural networks, great achievements have been made on 2D object detection, while 3D object detection is still underdeveloped. For example, according to the KITTI Object Detection Benchmark[43], the Average Precision (AP) of top 10 2D car detection algorithms is over 90% whereas best performance of 3D car detection is only 73.66%. This gap results from the difficulty of adding the third dimension and the orientation of the 3D bounding box.

For AVs, 3D information of the surrounding vehicles is indispensable because it expresses the vehicle dimensions, locations and orientation in the real 3D environment, which is essential for AVs to perform planning and decision making. To find a safe and efficient route, the information of actual and potential movement, dimensions, and location of other vehicles is necessary. In order to perceive the movement, 3D localization, orientation, and time are used to recover the velocity. The decision-making systems are more complicated and require more detailed 3D information, for example, the Bosch's Autonomous Emergency Braking systems (AEB) requires the distance to each part of a foregoing vehicle to decide whether or which level of brakes to apply. And it is common that some parts of the vehicle are occluded by other objects or truncated by the boundaries of the image. Thus, the exact location of each vehicle part and the visibility property of these parts are necessary.

Here we propose an approach that can simultaneously present the 3D dimensions, 3D localization, 3D bounding box orientation, 3D parts location, and parts visibility of a vehicle by giving a monocular image and all the 2D bounding boxes for vehicles. Figure 1 shows one output image example. TBD

1.2 Contributions

The first contribution is that our approach can predict not only the 3D bounding box for each vehicle but also the 3D position of each vehicle part, even if these parts are occluded by other objects or truncated by the boundaries of the image.

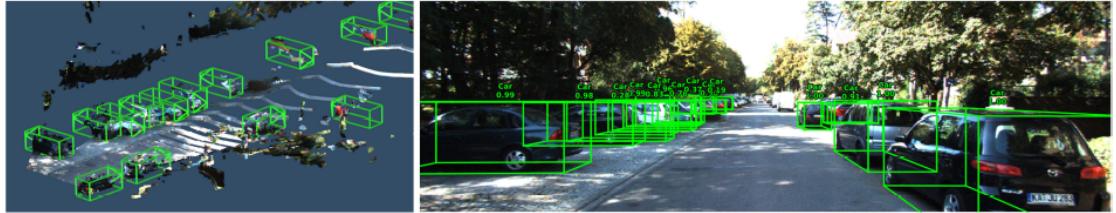


Figure 1: a

The foundation behind this is that vehicles are rigid objects and their geometric characteristics have a lot in common despite of the vehicle types. Therefore, these shared geometric characteristics serve as the priors which make it reasonable that each vehicle can be expressed by a 3D artificial model along with a scaling vector. This 3D model is integrated by vehicle parts which are further encoded by characteristic points. Besides, we apply regression rather than detection to search and locate these characteristic points. By doing so, our approach can find all the rest of parts, as long as it ascertains the existence of the vehicle object and some parts. Therefore, even through some parts are invisible to the camera, our approach can still localize them.

The second contribution is the proposed semi-automatic labelling process which generates additional labels to create a new dataset for our approach. Deep neural networks are greedy for data. Manual labelling is time-consuming and error-prone, let alone it is infeasible to correctly label the very far vehicles or occluded parts in the image. Therefore, we eliminate the human labor to the extend where only the 3D models need to be labelled manually. Then the process automatically selects the best-matching model and projects this model into the real image to generate labels, *e.g.* visibility and characteristic points. Besides, this process can be easily generalized to other tasks, as long as they require 2D geometry information in the image coordinate and 3D geometry knowledge in the world coordinate.

The third contribution is the multi-task framework which includes a prediction neural network and a inference block. The network can simultaneously perform vehicle parts localization, visibility characterization, and template proximity prediction at high accuracy level and consequentially, the inference block performs the vehicle dimension estimation, 3D vehicle localization, and rotation recovery. All these tasks can be completed in real time, 0.02s per image, which makes it applicable to directly process the video sequence of the vehicle's front camera.

2 Related Work

[what written in related work is 3D object detection, is it compatible with 3D pose estimation. the name of 3D vehicle detection in Exp.] TBD

2.1 Mono RGB Image Based Approaches

Most vehicles are equipped with monocular cameras and RGB images have detailed texture information with high resolution so that plenty of approaches are developed based on monocular RGB images. Images are mainly contain 2D information so most approaches in this category require other helps.

One way is to utilize the geometry information of the vehicles. In this category, a set of algorithms recover the 3D bounding box based on 3D car models. [120, 121] model the 3D geometry representation of objects with 3D wireframe models. And the 3D bounding box is estimated based on the geometry constraints of the vehicle in image and its corresponding 3D wireframe model. 3DVP [112] performs 3D detection based on 3D Voxel Patterns which are generated from the KITTI dataset [43] and a set of 3D CAD models to encode the geometric information of objects. Mono3D [26] introduces 3D proposals. It exhaustively places 3D bounding boxes on the ground-plane as proposals, then scores each proposal based on several hand-crafted geometric features, and applies a CNN to score the most promising candidates to generate the final 3D detections. Recently, CNNs are introduced to estimate some key features for 3D detection. Deep3DBox [74] applies two CNNs to estimate the orientation and dimensions respectively, makes use of the geometric criterion that a 3D box should fit tightly within the 2D box of the vehicle to estimate the translation, and finally integrates them together to perform 3D detection. Deep MANTA [24] generates 3D detection based on 2D detection. It applies a CNN to predict some 2D key points and the template similarity of the vehicles in image and recovers the corresponding 3D key points and 3D dimensions with some 3D CAD vehicle models, and finally perform 3D detection via 2D-3D matching. Our approach is similar to Deep MANTA in spirit.

Another way is to use temporal information. [34, 97] make use of motion structure and ground estimation to perform 3D detection from 2D bounding box.

2.2 Stereo RGB Image Based Approaches

A pair of stereo RGB images can provide the depth information of the current scene. 3DOP [28] recovers depth from stereo images and generates 3D box proposals with the constraints from depth and other geometry characters, which are forwarded to modified Fast R-CNN [44] pipeline for object detection and pose estimation. [82] extends 3DOP by applying a separate CNN to extract features from depth and integrating them together for the final 3D detection.

2.3 Point Cloud Based Approaches

A LiDAR point cloud is a collection of 3D points acquired by a LiDAR laser scanner. It can represent 3D information of the surroundings but its resolution is lower than images’.

One set of approaches apply a point cloud by converting it into a 2D array and making use of the image-based methods to perform 3D detection tasks. This can alleviate the inherent problem that LiDAR point is often sparse and irregular in 3D voxel representation . VeloFCN [68] projects the point cloud to the front view to generate a 2D point map and then applies a fully convolutional neural network to estimate the 3D bounding boxes for vehicles on this 2D point map. [111, 27] also fall into this category. Besides, the pre-trained models based on RGB images can be used to initialize the CNNs, which has been proven to be beneficial [49]. However, the 2D representation of LiDAR cloud points suffers from object size variations because of its distance to the sensor and object overlapping.

Another category of approaches transforms the cloud point into a 3D voxel grid representation. Various quantities are used to encode voxels [99, 98, 67, 105, 37], For example, in [105, 37], the information in one non-empty cell is encoded with six quantities while empty cells contain no information. For the single-stage detectors, *e.g.* Sliding Shapes [99] and Vote3D [105], slide window across the 3D voxel space and apply SVM classifiers to perform 3D detection. To improve the performance, Vote3Deep [37] uses a voting strategy with sparse 3D convolution. And [67] feeds the point cloud voxel directly to a 3D FCN to generate 3D bounding boxes. For two-stage framework, *e.g.* Voxelnet, extends the 2D RPN [87] to 3D to generate 3D proposals and apply a refine network to score these proposals [119]. 3D CNNs boosts the 3D detection performance while the computation is very expensive and sparsity is to be explored.

2.4 Sensor Fusion Based Approaches

Images are of high resolution but lack depth, while LiDAR data has 3D information but are sparse and irregular. Therefore many works have investigated combining these complementary data sources together to build robust 3D object detectors.

One investigated direction is to extract information from these two sources separately for sub-functions. Frustum PointNets [84] first generates frustum point cloud proposals based on 2D object detection, then performs 3D instance segmentation in each frustum point cloud and finally applies a box estimation net to estimate the final 3D bounding box for each object. Similarly, [35] first generates 2D bounding boxes and estimate vehicle dimensions for the target vehicles on images. Secondly, a model fitting algorithm estimates the 3D bounding box on a subsets of the point clouds which fall into the 2D bounding box after projection. Finally a refine 2D CNN performs the final 3D box regression and classification.

The other is to fuse these two kinds of data and process it as a whole. There are various ways to fuse data. One intuitive way is projection, *e.g.* [36] converts the point cloud to a dense depth image and then appends it as an additional channel of image, [91] extends [36] via converting the cloud point to a three-channel HHA map. Their 3D detection score is low mostly due to the lose of 3D information during projection. The more advanced way is to fuse their feature maps. MV3D [29] first generates 3D proposals in LiDAR’s bird’s eye view and projects these proposals to the features maps of the image and the bird’s eye view and front view of LIDAR data. Then a deep network fuses three ROI pooling regions and performs 3D detection based on the fused features. PointFusion[114] applies a ResNet [52] to extract appearance and geometry features from image crops defined by the 2D bounding boxes and a modified PointNet [85] to process the raw point cloud simultaneously, and finally uses a novel fusion network to integrate both features and estimate 3D bounding box. Recently, Zining *et al.* [108] applies an innovative sparse non-homogeneous pooling layer to fuse the features extracted from bird’s eye view of LiDAR data and front view camera images by two separate CNNs before region proposal stage. And then a single-stage detector adapted from [71] is designed to perform 3D detection on these fused data without RoI pooling, which is 6 times faster then MV3D. AVOD [64] applies double feature fusions. It uses two identical CNN to extract features from images and LiDAR data respectively and exploits a Multimodal Fusion Region Proposal Network to generate vehicle proposals from the fused feature crops, projection regions of anchors in the feature maps. Finally it applies a Multiview Detection Network to perform 3D detection on the fused feature crops corresponding to the proposals. It achieves the best

performance on KITTI 3D vehicle detection up till now.

3 Background

3.1 A Short History of Autonomous Vehicles

Autonomous Vehicles (a.k.a. Automated/Driverless/Unmanned/Robotic/Self-driving Vehicles) is a relatively vague concept to the public. Actually, it covers a continuum from traditional fully human-driving automobiles to fully self-driving vehicles, as SAE classified in the Table 1.

SAE level	SAE name	SAE narrative definition	Execution of steering and acceleration/deceleration	Monitoring of driving environment	Fallback performance of dynamic driving task	System capability (driving modes)	BAsT level	NHTSA level
<i>Human driver monitors the driving environment</i>								
0	No Automation	the full-time performance by the <i>human driver</i> of all aspects of the <i>dynamic driving task</i> , even when enhanced by warning or intervention systems	Human driver	Human driver	Human driver	n/a	Driver only	0
1	Driver Assistance	the <i>driving mode-specific execution</i> by a driver assistance system of either steering or acceleration/deceleration using information about the driving environment and with the expectation that the <i>human driver</i> perform all remaining aspects of the <i>dynamic driving task</i>	Human driver and system	Human driver	Human driver	Some driving modes	Assisted	1
2	Partial Automation	the <i>driving mode-specific execution</i> by one or more driver assistance systems of both steering and acceleration/deceleration using information about the driving environment and with the expectation that the <i>human driver</i> perform all remaining aspects of the <i>dynamic driving task</i>	System	Human driver	Human driver	Some driving modes	Partially automated	2
<i>Automated driving system (“system”) monitors the driving environment</i>								
3	Conditional Automation	the <i>driving mode-specific performance</i> by an <i>automated driving system</i> of all aspects of the <i>dynamic driving task</i> with the expectation that the <i>human driver</i> will respond appropriately to a <i>request to intervene</i>	System	System	Human driver	Some driving modes	Highly automated	3
4	High Automation	the <i>driving mode-specific performance</i> by an <i>automated driving system</i> of all aspects of the <i>dynamic driving task</i> , even if a <i>human driver</i> does not respond appropriately to a <i>request to intervene</i>	System	System	System	Some driving modes	Fully automated	3/4
5	Full Automation	the full-time performance by an <i>automated driving system</i> of all aspects of the <i>dynamic driving task</i> under all roadway and environmental conditions that can be managed by a <i>human driver</i>	System	System	System	All driving modes		

Table 1: Summary of levels of driving automation[9]

Contrary to the intuition, the idea of Autonomous Vehicles has a long history. The experiments of this fictional idea can be tracked back to 1920s and the technology behind it was radio control [8].

But the truly autonomous cars did not show up until 1980s, even though they could only move slowly on clear streets and required large human intervention.

Mercedes-Benz demonstrated a robotic van based on saccadic vision [92]. The Autonomous Land Vehicle (ALV) project funded by The Defense Advanced Research Projects Agency (DARPA) cultivated automatic vehicles based on Computer Vision, LIDAR, and autonomous robotic control [62]. Carnegie Mellon University initiatively applied neural network to control the vehicle [83].

In 1990s, huge progress was made. The VaMP from Daimler-Benz drove more than 1000 km, achieving the maximum speed of 130 km/h on a normal Pairs highway semi-autonomously [92]. The Navlab project in Carnegie Mellon University achieved a 5000 km journey cross America with only 1.8% human interventions [5]. The ParkShuttle in Netherlands could autonomously navigate itself on a dedicated lane as an automated people mover [79]. In this decade, the experiments were mainly carried out in highway scenarios rather than urban scenes.

In 2000s, competitions promoted this technology a lot. One of the most famous is The DARPA Grand Challenge in U.S. who offered \$ 1 million for the first prize. In 2004, no vehicle completed the 241-km journey autonomously while 5 teams achieved this goal in 2005 [60]. And in Grand Challenge III 2007 , known as Urban Challenge, 6 vehicles finished the event which was a 96 km urban route involving with traffic regulations and other vehicles [61].

In 2010s, Autonomous Driving Technology starts to take off. Numerous events and projects have been carried out and considerable companies, universities, and research centres have engaged into this field. Based on the progress made before, many autonomous vehicle systems are being tested or even brought into production. Notable events covers the VisLab Intercontinental Autonomous Challenge in 2010 [21] and the Intelligent Vehicle Future Challenges from 2009 to 2013 [75]. In industry, Tesla Motor released AutoPilot that is able to perform automated parking and lane control with autonomous driving, braking and speed adjustment in 2014, and Audi started the first production car, A8, reaching Level 3 of Automation in 2017 [4].

Based of the efforts in the last decades, Autonomous Driving gradually transform from dream to reality. Its bonus covers safety guarantee, congestion reduction, land use efficiency, energy and emission reduction, economic benefits and so on. However, the Level 5 vehicles are so far from maturity that the research and development are in high demand.

3.2 Deep Learning Technology

Recently, Deep Learning (DL) has shown its impressive power in a variety of tasks, especially in some complex tasks that can not be explicitly programmed by hand,

such as anomaly detection and online advertising. DL delivers the solutions by learning from data automatically via a general learning procedure which dominantly makes use of backpropagation and optimization algorithms, *e.g.* gradient descent.

It attracts the world-wide attention mainly by outperforming other classic machine learning algorithms, *e.g.* Support Vector Machine (SVM), in many competitions involved with image classification, object detection, or nature language processing. This is because DL applies a deep automatic feature learning architecture, *i.e.*, a artificial neural network model with multiple hidden layers, to learn deep distributed hierarchical nonlinear representations which yield better performance for learning tasks, *e.g.* in terms of classification accuracy. Such representations bring many good properties, such as feature reuse, parameter sharing, multiple levels of progressive abstraction, and invariance to local changes of the raw inputs, and thus provide better predictive power than classic machine learning algorithms [15].

Convolutional neural networks (CNNs) is one of the major branches of artificial neural networks. A CNN is a feed-forward multi-layer neural network, typically consisting of one or more convolutional layers, interleaved by some pooling layers, and finally followed by some fully-connected layers. This modern framework of CNNs is established by LeCun *et al.* when he proposed the handwritten digit classifier LeNet-5 [65]. Afterwards, more and more deeper architectures are explored such as AlexNet [90], VGGNet [93], GoogleNet [101] and ResNet [53]. In general, deeper architectures, bring better feature representations, and closer approximations to the target function, as well as more complex models which are more difficult to train and easier to be overfitting. Therefore many approaches have been proposed to address such problems.

CNNs are specifically designed to work with problems taking images as inputs, such as image classification, object detection and pose estimation. The above-mentioned CNNs examples all made great performance in their respective tasks, *e.g.* ResNet won the first prize in ILSVRC 2015.

The core building block, convolutional layers, is used to compute feature maps with convolution kernels. Each neuron in a convolutional layer is connected to a local region in the previous layer, called receptive field, and computes an output by performing an convolution operation (element-wise matrix multiplication) between its weights (kernel) and the connected region followed by an non-linear activation function. A nice property of CNNs is that the kernel is shared by all receptive fields in the preceding layer when computing the corresponding feature map. Thus, each feature map is used to capture exactly the same feature but at

different locations. And this characteristic can substantially reduce the number of parameters in CNNs, which will lead to faster training [77].

The pooling layers perform a downsampling operation, typically max pooling [20] and average pooling [106], along the spatial dimensions of feature maps to achieve shift-invariance.

Normally, successive convolution layers detect more abstract features, *e.g.* wheels, than the preceding ones that tend to detect low-level features, *e.g.* curves and edges. Therefore, after the operations performed by several convolutional and pooling layers, progressively higher-level features can be obtained to feed a fully-connected layer whose goal is to perform high-level reasoning to generate the global semantic information [93, 54].

For classification problems, the output layer of CNNs usually deploys the softmax operator [90] or the SVM method [104] to output discrete results. While regression tasks require continuous-valued predictions so that the output layer should have a linear activation function, *e.g.* weighted sum, along with a proper cost function, *e.g.* mean squared error [118].

The training for CNNs is a global optimization problem by minimizing the defined loss function. Normally, CNNs can be trained end-to-end efficiently with back-propagation together with an optimization algorithm, such as stochastic gradient descent [109]. The mechanism behind it is that gradient of the loss function w.r.t. all parameters is calculated and then used to update the parameters to the direction of minimizing the loss function based on iterations over the full batch or mini batches of the training dataset.

3.3 Computer Vision

Computer vision is an interdisciplinary field that enables computers to interpret images just as what we human can do. Its goal is about automatic extraction, analysis, and understanding of the information from images [7]. There are a huge variety of ways to process images and a marvellous diversity of applications in this board field, ranging from replicating human visual abilities to creating nonhuman visual capabilities [48]. Some real-world applications include object recognition or detection, 3D model building, motion capture, etc [103].

Despite all the progress achieved, computer vision systems are still so underdeveloped that they can't match the visual ability of a 5 year-old child. Vision is natural and effortless for humans but factitious and demanding for computers [12]. The difficulty in part results from that vision is a inverse problem where we at-

tempt to recover the unknowns, *e.g.* shape and illumination, based on insufficient information of their causes, *e.g.* models [103]. Therefore, it is hard to conclude the clear rules that can help get a solution explicitly, especially for complex problems, *e.g.* 3D object detection. But currently, Deep Learning seems to find a way out and pushes a big step forward.

Next in this section, we introduce some basic knowledge in computer vision related to our approach.

3.3.1 Image Formation

Imaging systems or cameras are some devices that allow the projection of light from 3D points to some 2D medium that records the light pattern. Figure 2 (a) shows a pinhole imaging model which is able to capture the light pattern, the inverted candle image, by allowing a very tiny cone of rays issued at every point of the source, the real candle, to project on the image plane. This is called pinhole perspective projection which is primitive but provides an acceptable approximation of the imaging process [42].

The projection equation can be derived from Figure 2 (b), where (O, i, j, k) is the pinhole camera coordinate system and the origin O is at the pinhole, $p = (u, v, d)^T$ is a point in image, and $P = (x, y, z)^T$ is a point denotes the source. As light travels straight in the same medium, p, O and P are collinear, and therefore we can deduce

$$\begin{cases} u = d \frac{x}{z}, \\ v = d \frac{y}{z} \end{cases} \quad (1)$$

Modern cameras are built on lenses that can gather more light to make the image more bright while maintaining its sharpness. But the imaging process is very similar to the pinhole camera. Lenses can also introduce some aberrations, *e.g.* spherical aberration, radial distortion, and chromatic aberration [42]. Therefore, correction is necessary.

3.3.2 Intrinsic and Extrinsic Parameters

To project a 3D point in the world coordinate system, first we have to transform this point to the camera coordinate system and then transform it to the image plane. The first transformation depends on the extrinsic parameters while the second rely on intrinsic parameters.

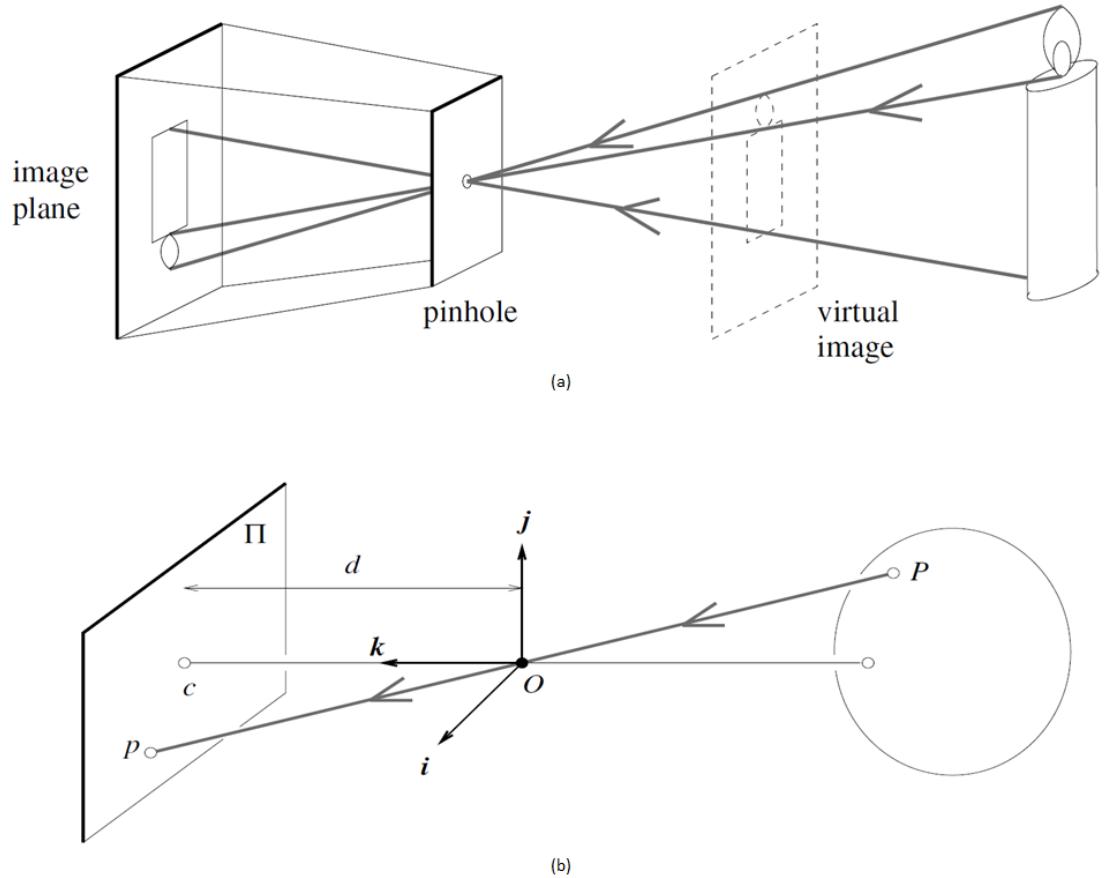


Figure 2: (a). the pinhole imaging model; (b). geometric model for perspective projection. [42]

Intrinsic parameters include the focal length f just like the d in Figure 2 (b), the image coordinates origin (u_0, v_0) , and the skewed angle θ of two image axes. Coordinates in image plane are usually expressed in pixel which can be square or rectangular. So let's assume α and β are the value expressed f with horizontal and vertical pixel-meter scales. Therefore, we can finally transform Eq. 1 into

$$\begin{cases} u = \alpha \frac{x}{z} - \alpha \cot \theta + u_0, \\ v = \frac{\beta}{\sin \theta} \frac{y}{z} + v_0 \end{cases} \quad (2)$$

When it is written in matrix form as Eq 3, the 3×3 matrix K is the intrinsic matrix of the camera. Note that P is expressed in camera coordinate system and p is in image coordinate system.

$$p = KP \Leftrightarrow \frac{1}{z} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} \alpha & -\alpha \cot \theta & x_0 \\ 0 & \frac{\beta}{\sin \theta} & y_0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} \quad (3)$$

Extrinsic parameters defines a rigid transformation from world coordinate system to camera frame. A rigid transformation has six degree of freedom, including three Euler angles expressed in a 3 rotation matrix R and three translation components along each axis expressed in a 1×3 translation vector t . A point expressed in homogeneous coordinates is $P = (x, y, z, 1)^T$. Homogeneous coordinates can simplify various geometric transformations into matrix multiplication [42]. Therefore, this rigid transformation expressed in homogeneous coordinates is

$$P^C = T_W^C P^W, \text{ where } T_W^C = \begin{pmatrix} R & t \\ 0^T & 1 \end{pmatrix} \quad (4)$$

P^C and P^W are coordinates of the same point expressed in world coordinate system and camera coordinate system respectively. T_W^C is the extrinsic matrix of the camera.

To put it all together, the projection equation in homogeneous coordinates is

$$p = \frac{1}{z} M P^W = \frac{1}{z} K \begin{pmatrix} R & t \\ 0^T & 1 \end{pmatrix} P^W \quad (5)$$

where M is the **perspective projection matrix**.

4 My Work - Technical Detail

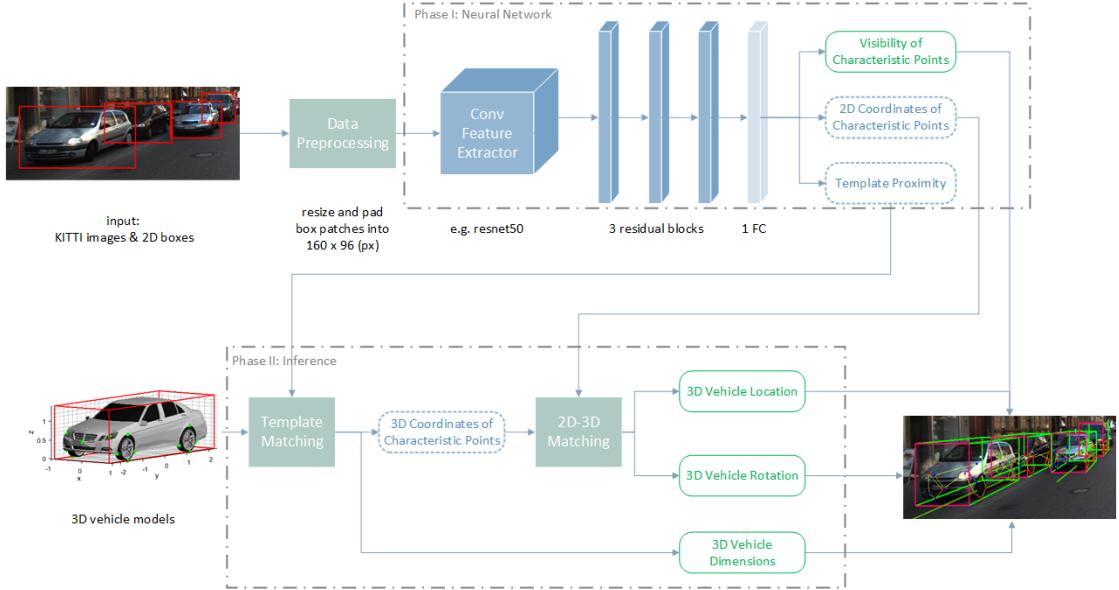


Figure 3: The Architecture of our Approach

In this section, we elaborate our approach for 3D pose estimation of vehicles in monocular images. Our approach is based on 2D vehicle detection, *i.e.*, the 2D bounding boxes of vehicles are given as prior, since 2D object detection techniques, such as YOLO [86], Faster R-CNN [87], SSD [72], R-FCN [32], and FPN [70], have achieved very high accuracy while end-to-end 3D object detection performance is still primitive. The overall architecture of the approach is illustrated in Figure 3. It consists of two main phases. The first one is the CVT Network which takes the KITTI images and the 2D bounding box as inputs and predicts the visibility property and 2D coordinates of all characteristic points, as well as the template proximity. The details are presented in section 4.2. The other phase estimates the 3D location, rotation, and dimension of each vehicle based on the 3D vehicle dataset and the outputs of the CVT Network. Section 4.3 demonstrated this phase in detail. The KITTI dataset and the 3D vehicle dataset are described in section 4.1.

4.1 Data and Labelling

4.1.1 KITTI Dataset

The KITTI 2D/3D object detection challenge is dedicated to autonomous driving and releases a dataset containing 7481 images for each one of 4 cameras and their associated labels and calibrations [43]. The dataset covers scenarios of city, residential, road, campus, and person. We only use the images taken by the left colour camera. The example image is shown in Figure 4. The associated labels are show in Table 2. The calibration is given as transformation matrix.



Figure 4: Example image from KITTI dataset captured by the left colour camera

#Values	Name	Description
1	type	Type of object: Car, Van, Truck, Cyclist, Tram, Pedestrian, Person_sitting, Misc or DontCare
1	truncated	Float from 0 (non-truncated) to 1, indicating the extent of the object out of the image
1	occluded	Integer indicating occlusion state: 0 = fully visible, 1 = partly occluded, 2 = largely occluded, 3 = unknown
1	alpha	Observation angle of object, ranging [-pi..pi]
4	bbox	2D bounding box of object in the image (0-based index): contains left, top, right, bottom pixel coordinates (x_1, y_1, x_2, y_2)
3	dimensions	3D object dimensions: height, width, length (in meters)
3	location	3D object location x,y,z in camera coordinates (in meters)
1	rotation_y	Rotation ry around Y-axis in camera coordinates [-pi..pi]

Table 2: Label specification of KITTI dataset

4.1.2 3D Vehicle Dataset

This dataset is intended to encode the diversity of vehicles by means of type, dimension, and chassis shape. It is created based on a selected subset of synthetic 3D CAD models[40]. We only include the well-aligned and symmetrical models, and also take the dimension and type into consideration. We classify all the models into six categories: Mini, Hatchback, Sedan, SUV, Wagon, and Van. And for each category, subcategories are marked out according to their dimensions and shapes in order to make a refined selection. The statistic of dimensions is well distributed and each existing car can be classified into one category exclusively.

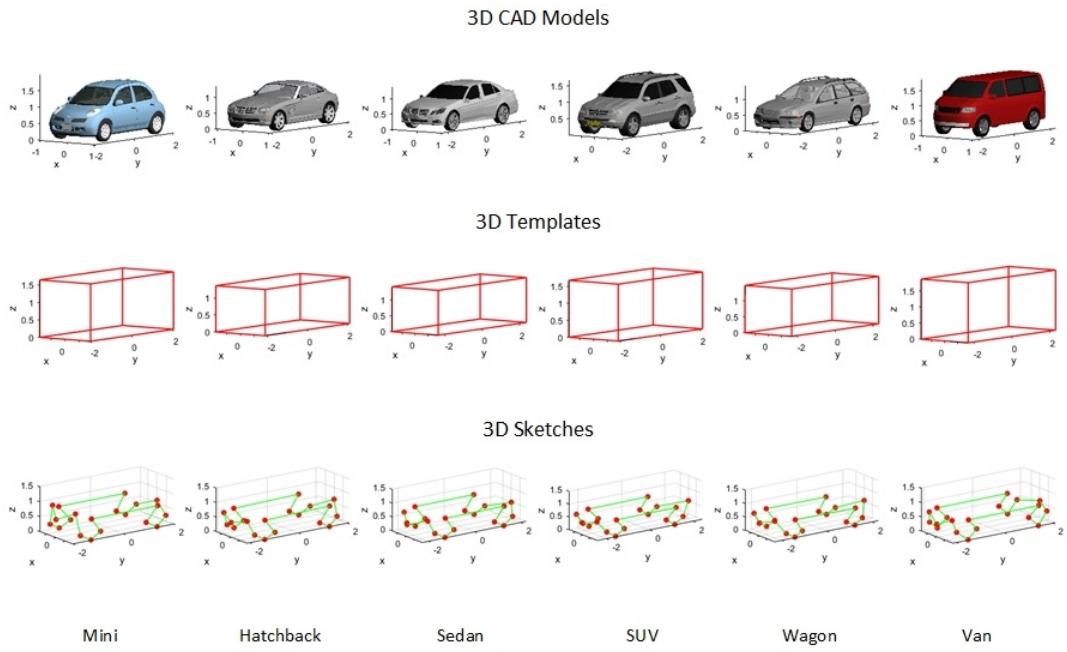


Figure 5: Six vehicle categories of the 3D vehicle dataset. Each vehicle model is associated to a 3D CAD model, a 3D template, and a 3D sketch.

Our final dataset consists of 54 valid vehicle models. Each vehicle model has a 3D CAD model, a 3D template, and a 3D sketch correspondingly, as shown in Figure 5. All these three representation are aligned in canonical view and share an identical object coordinate system. The CAD models are created based on real vehicles and have all the geometry information with them. The 3D templates represent the vehicles' dimensions. The 3D template associated to the 3D model k is denoted as $t_k = (h_k, w_k, l_k)$ where h_k , w_k and l_k represent the height, width, and length of the model respectively.

The 3D sketches indicate the chassis shapes. Each 3D sketch consists of 20 charac-

teristic points around the chassis with each point denoting one part of the vehicle. So the k^{th} sketch is denoted as $S_k^{3d} = (p_1, p_2, \dots, p_{20})$, where $p_i = (x_i, y_i, z_i)$. The reason why we choose the points around chassis as feature points is that their geometric relationship is more stable than points in other places. Because the chassis part is more about functionality than appearance attraction compared to other parts of the vehicle, *e.g.* the upper part. In this way, the sketches can be more universal so that they can represent a wider range of vehicles, which reduces the number of models required and further increases the computational efficiency.

In order to create the 3D sketches, we implement a labelling tool, shown in Figure 6. It is capable to label the 3D coordinates for all the characteristic points associated to each vehicle. [should it be elaborated in detail]

TBD

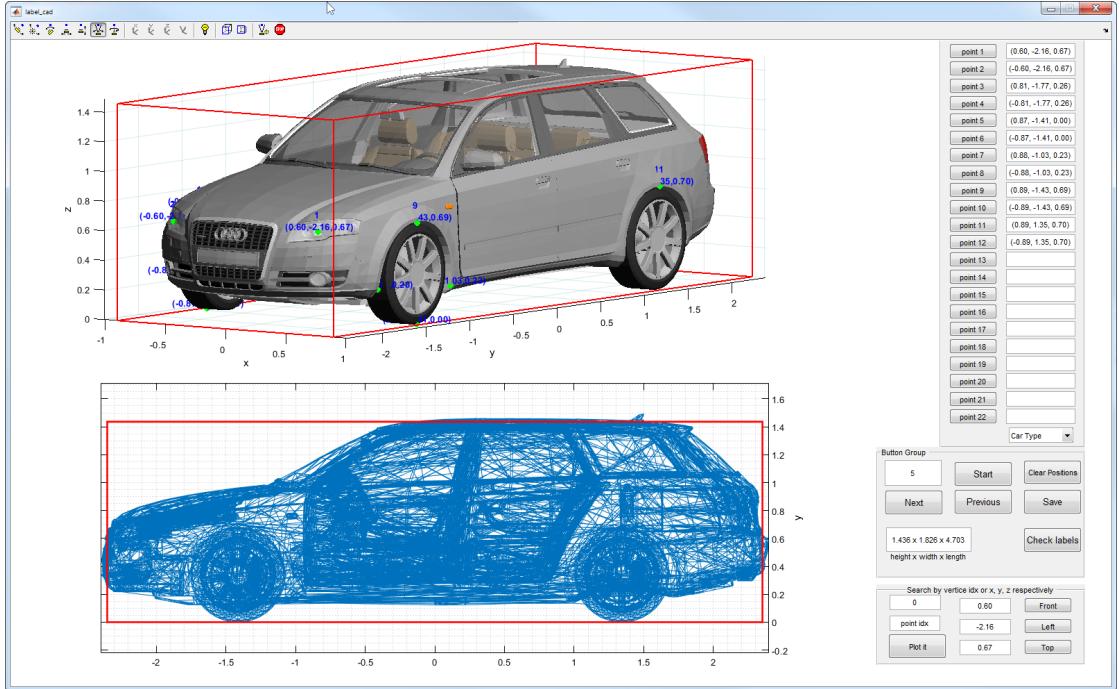


Figure 6: Labelling tool for 3D sketches

4.1.3 Labelling for KITTI Dataset

Our approach requires three extra kinds of labels to training the CVT Network, *i.e.*, 2D coordinates and visibility property of characteristic points and template proximity, as shown in Table 3. Labelling is never a trivial task. Due to the demanding workload of manual labelling and the cases where it is almost impossible

to label the small vehicles in image manually, we propose an automatic label generation method. It is able to make use of the 3D vehicle dataset and the KITTI dataset to generate these three additional kinds of ground truth.

#Values	Name	Description
2x20	2D coordinate	(x, y) , location of 20 characteristic points in image coordinates
1x20	visibility	Integer indicating visibility property of each point: 0 = visible, 1 = occluded, 2 = self-occluded, 3 = truncated
3x54	template proximity	A vector T_i represents the dimension ratios between each model and the vehicle

Table 3: Specification of three additional labels

4.1.3.1 2D Coordinates

2D coordinates of key points of each vehicle in image coordinates is generated by projecting the 3D sketch of this vehicle to the image. The vehicle’s 3D sketch is selected via template-matching, *i.e.*, the best-matching sketch is the one whose associated template is closest to the vehicle’s dimensions. The 3D-2D projection is performed based on the given intrinsic and extrinsic parameters, as shown in Figure 7. The intrinsic parameters are given as calibration by KITTI and the extrinsic parameters are given as 3D object dimensions and rotation ry in KITTI labels. KITTI makes a simplified assumption here that the vehicle only rotates around the yaw axis but not roll or pitch axis. One example of the labelled 2D coordinates is shown in Figure 8.

4.1.3.2 Visibility

As an example of the labelled visibility shows in Figure 8, the visibility property of characteristic points is classified into four scenarios:

- i. Visible if the point can be seen directly;
- ii. Occluded if the point is occluded by other objects;
- iii. Self-occluded if the point is blocked out by the vehicle itself;
- iv. Truncated if the point exceeds the boundaries of the image.

The visibility of each point is determined by its position, *i.e.*, the 2D coordinate in image. Being visible or self-occluded is distinguished by means of rotation ry of each vehicle. As the right plan sketch in Figure 9 shows, the rotation ry can

$$\begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{bmatrix} f & s & x_c \\ 0 & af & y_c \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} R_{3 \times 3} & 0_{3 \times 1} \\ 0_{1 \times 3} & 1 \end{bmatrix} \begin{bmatrix} I_{3 \times 3} & T_{3 \times 1} \\ 0_{3 \times 1} & 1 \end{bmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$

Points
 in 2D image
 coordinates Intrinsic Projection
 Rotation Translation
 Points
 in 3D object
 coordinates

Transform from camera
 coordinates to image
 coordinates Transform from object
 coordinates to camera
 coordinates

Figure 7: 3D-2D projection

indicate unambiguously which faces of the vehicle are facing to or away from the camera. If the points are on the observable faces, they are labelled as visible, otherwise they are classified as self-occluded. Occluded case happens when the 2D coordinate of the point falls into the 2D bounding box of a former object. The object is defined as former when it locates in the region enclosed by the axes and two solid blue line in the left schematic plot of Figure 9 if the vehicle is on the first quartile. And when the vehicle is on the second quartile, it's just a mirrored case. Truncated is defined that its 2D coordinate exceeds the boundaries of the image. The size of KITTI images are determined which makes it easy to implement. The truncated property has the highest priority, the occluded underlies , the self-occluded or the visible is considered at last.

4.1.3.3 Template Proximity

Template proximity of one vehicle is encoded as a vector T_i which is defined as $T_i = \{r_k\}_{k \in \{1, \dots, K\}}$, where K denotes the number of 3D vehicle models and $r_k = (r_h, r_w, r_l)$ corresponds to three scaling ratios between the dimensions (*i.e.*, height, width, and length) of each model and the vehicle's respectively. The vector T_i represents the similarity between each model and the vehicle. The most similar model is the one whose r_k is closest to $(1, 1, 1)$.



Figure 8: Example of 2D coordinates and visibility of one vehicle. The left image is one patch of KITTI image. The right image is the one after labelling. The points indicate the 2D coordinates and the color indicates the types of visibility: red for visible, green for occluded, and blue for self-occluded.

4.1.4 Notations

In sum, based on the KITTI image and 3D vehicle models, each vehicle can be defined by seven critical attributes:

$$\{D, B^{2d}, B^{3d}, C^{2d}, C^{3d}, V, T\}$$

$D = (h, w, l)$ represents the dimensions of the vehicle. $B^{2d} = (c_u, c_v, w, h)$ defines the 2D bounding box in image with (c_u, c_v) denoting its center, w for its width and h for its height. $B^{3d} = (o, \theta, d)$ where $o = (c_x, c_y, c_z)$ is the center, θ is the rotation ry around the yaw axis, and $d = (w, h, l)$ is its dimensions. $C^{2d} = \{(u_i, v_i)\}_{i \in \{1, \dots, 20\}}$ represents the 2D part coordinates in image plane, while $C^{3d} = \{(x_i, y_i, z_i)\}_{i \in \{1, \dots, 20\}}$ denotes the 3D part coordinates in world coordinate system. $V = \{v_i\}_{i \in \{1, \dots, 20\}}$ is the visibility for all the characteristic points in the vehicle. $T = \{(r_h, r_w, r_l)_k\}_{k \in \{1, \dots, 54\}}$ is the template proximity vector.

4.2 Deep CVT Network

This subsection describes the details of the Deep CVT Network, including its architecture, implementation, and training.

4.2.1 Architecture

The neural network is to learn a map from an N-dimensional input space to an M-dimension output space. The map consists of several stages, called layers of the

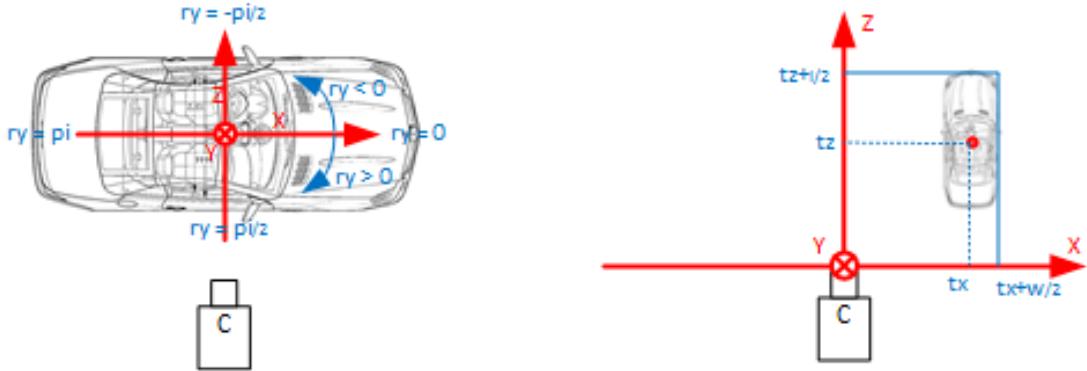


Figure 9: Visibility labelling mechanism. The left image shows the rotation ry in the camera coordinate system, which helps distinguish visible or self-occluded case. The left image shows the location of the vehicle in the camera coordinate system, which helps defined the occluded situation.

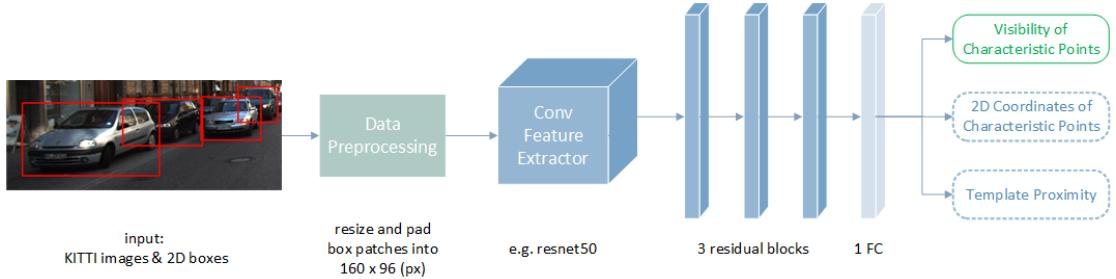


Figure 10: The Architecture of Deep CVT Network

network, written as

$$Y = l_L(\dots l_2(l_1(X))), \text{ where } Y \in \mathbb{R}^M, X \in \mathbb{R}^N \quad (6)$$

Even though there are various layers, most layers is composed of neurons, the basic computation element and most neurons are realized by linear and non-linear operations as

$$a_{ij} = \sigma(W_{ij}^T X_i + b_j) \quad (7)$$

where a_{ij} is the result of the i th neuron in the j th layer, σ is the activation function and W_{ij} and b_j are the weights vector and bias applied to input vector X_i . The pooling layer is a special case which performs downsampling operations, such as max pooling [20] and average pooling [106].

As shown in Figure 10, our network follows the standard CNN architecture established in LeNet-5 [65]. It stacks a set of convolutional layers and pooling layers as

a feature extractor, followed by 3 residual blocks and one fully-connected layers TBD to refine the features, and three kinds of output layers in the end. The goal of this network is to learn the 2D coordinates and visibility property of 20 interest points and the template proximity of the vehicle, given the RGB images captured by a monocular camera.

4.2.1.1 Input Layer

The input layer accepts the input RGB images and feeds them into the network. In theory, images with arbitrary size can be accepted, but to make it more efficient, we use images with fixed resolution (96 x 160). In this way, multiple images can be processed in one batch in order to reduce the variance of parameter updates and make the best use of highly optimized matrix optimizations during training [88].

Since our approach is based on 2D object detection, each image contains at least one whole vehicle . To obtain these images, we first crop the patches defined by the 2D bounding box, then resize them to match one predefined dimension (96 or 160), and finally put the resized images in the center of the 96 x 160 canvas and padded with zeros for other pixel positions. We choose 96 x 160 as the fixed size because they are the means of sizes of all original patches so that we don't have to rescale the patch too much. Some examples are shown in Figure 11.



Figure 11: Examples of input images: a. the side view, b. the back view, c. the front view, d. the occluded case

4.2.1.2 Hidden Layers

The hidden layers in our network consist of a feature extractor, 3 residual blocks [52], and a fully-connected layer. Both the feature extractor and residual blocks are composed of convolutional layers and pooling layers.

The feature extractor is actually one of available neural network models built in the Keras library [30], *i.e.*, VGG16, VGG19, Xception, ResNet50, *etc.* It is trained to learn deep distributed hierarchical nonlinear representations of the input images.

ResNet50 [52] is chosen to be the benchmark model because it can ease the training and accelerate the convergence, especially for the fine tuning. It is well known that there are two main obstacles that hamper the training of deep neural networks: the vanishing/exploding gradient problem [16, 46] and the degradation problem that as the network is built deeper, the accuracy gets saturated or even degrades sharply [51]. Residual nets are designed to address these two problems. Instead of learning the direct mapping from input to output, it learns a residual mapping first and then add the input to it, which just as Figure 12 (a) shows. The paper validates that it is easier to learn the residual mapping than the direct one and the gradient can always pass through along the shortcut connections. Figure 12 (b) and (c) shows two building blocks in our network. (b) is an identity block where the dimensions of input and output matches, while (c) is a linear projection block where the right branch performs 1×1 convolution to match the dimensions because the final addition is performed element-wisely.

The last three residual blocks and a fully-connected layer is designed to learn nonlinear combinations of the extracted features and then forward them to the output layers.

For all hidden layers, we all apply ReLU as the activation function. The function is

$$f(z) = \max\{0, z\} \quad \text{where } z = W^T X + b \quad (8)$$

This is a piecewise linear function consisting of two linear pieces, which makes the gradients through a rectified linear unit stay large and consistent whenever the unit is active. Papers [47, 76, 58] have validated that networks activated by ReLU can achieve much better performance than others. ReLU preserves many properties that make the model easy to optimize with gradient-based methods and generalize well, while introducing non-linearity into the model [48].

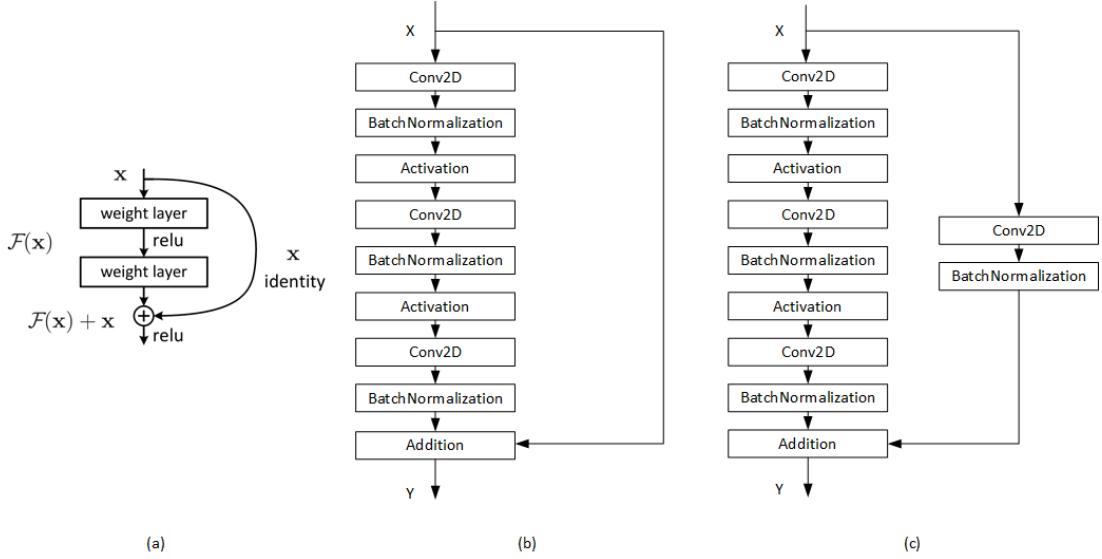


Figure 12: (a). a residual block [52], (b). an identity residual block in our network, (c). a linear projection residual block in our network

4.2.1.3 Output Layers

Our network has three kinds of output layer for the three tasks respectively. For 2D coordinates regression, a 40-neuron layer with linear activation function is provided. Each neuron's output indicates one coordinate value, *i.e.*, x_i or y_i . 20 coordinates are arranged in an ascending order. For visibility characterization, we implement 20 independent quaternary classifiers for 20 points. Each classifier is associated with a softmax activation function [19] which outputs four values representing the probabilities of each target class over all possible classes. Template proximity is equipped with a 154-neuron linear layer for regression. Each neuron indicates one ratio element of the template proximity vector. Eq.9 is a linear activation function while Eq.10 is a softmax activation function.

$$out_i = W_i^T X_i + b_i \quad (9)$$

$$out_i = \frac{e^{W_i^T X_i + b_i}}{\sum_{k=1}^K e^{W_k^T X_k + b_k}} \quad (10)$$

4.2.2 Training

4.2.2.1 Multi-task Learning

As mentioned in Section 4.2.1.3, our network has totally 22 output layers. This implementation makes use of multi-task learning (MTL) where multiple learning tasks are trained parallel based on the shared features. Caruana has validated that tasks trained in MTL have better generalization performance than trained in a single-task learning mode [22]. And he summarizes that this improvement results from leveraging the domain-specific information contained in the training signals of other related tasks [22]. Therefore, we follow this paradigm to design our network in order to achieve better performance.

4.2.2.2 Loss Functions

The learning of 2D coordinates and template proximity are regression tasks so that a robust smooth L_1 loss function [44] is chosen for them. We modify it as

$$\text{smooth}_{L_1}(x) = \begin{cases} x^2 & \text{if } |x| < 0.25 \\ |x| - 0.1875 & \text{otherwise} \end{cases} \quad \text{where } x = \hat{y} - y \quad (11)$$

Compared to L_2 loss, it is less sensitive to outliers and no special attention is required to pay in order to prevent gradient exploding problem[44]. Besides, when applied with gradient-based optimization, L_1 loss and L_2 loss often result in poor performance [48].

For visibility characterization, we develop a model for probabilistic classification so that categorical cross-entropy is used as the loss function. It is also known as the negative log-likelihood [48], defined as:

$$L_{cce} = \frac{1}{N} \sum_{n=1}^N H(y, \hat{y}) = -\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^c y_{n,i} \log \hat{y}_{n,i} \quad (12)$$

where $H(y, \hat{y})$ denotes the cross-entropy between the ground truth probability distribution y and the predicted \hat{y} , and $y_{n,i}$ represents the true probability of i_{th} class for the n_{th} data example $\hat{y}_{n,i}$ is the estimated. It is a continuous convex loss function which measures the discrepancy between the true and estimated distributions multi-class classification tasks [14] which means it can measure the degree of correctness, *i.e.*, it can distinguish between "nearly correct" and "totally

wrong” cases. Therefore, it outperforms other loss functions in classification tasks and then becomes ubiquitous in deep learning nowadays.

Therefore training objective is expressed as the total loss of all tasks, written as:

$$L_{total} = \lambda_{coord} L_{coord} + \lambda_{temp} L_{temp} + \lambda_{visib} \sum_{i=1}^{20} L_{visib\ i} \quad (13)$$

where λ_{coord} , λ_{temp} , and λ_{visib} are loss weights for these three kinds of tasks respectively. Loss with higher weights has more impact on the gradients and thus tunes the parameters perform better for its corresponding task.

4.2.2.3 Normalization

Normalization is an important pre-processing step in deep learning which ensures that each feature has a similar data distribution. This is usually done by restricting the features in certain range or standardizing their ranges. It can enhance the learning capability of the network and speed up the convergence substantially because it can reduce the bias among features and decrease the low and high frequency noisy in data [59]. The speed-up mechanism can be concluded from Figure 13. If we initialize the network at point A in Figure 13 (a), the gradient-descent update routine will oscillate along the long axis of the ellipse, which definitely takes more time to reach the optimum than any arbitrary initial point in Figure 13 (b) where the update trajectory is almost a straight line for any starting point to the minimum. Besides, this also helps the performance because the update trajectory oscillates less around the minimum for case (b) than case (a) so that it is more possible for case (b) to reach the optimum.

Our network takes RGB images as input so that pixel intensity is the input feature. The technique we use to normalize the image is channel mean subtraction, used in [94, 45], which centers all the features around the origin along each dimension. As mentioned in CS231N [3], mean channel subtraction is enough for CNNs and the original range of pixel values is determined, *i.e.*, [0, 255].

The range of ground truth influences the loss. In order to make the three labels impact similarly on the loss, we normalize them. For visibility, we use one-hot-encoding [6] to encode its classes so that it only has value 0 or 1. The value of 2D coordinates $C^{2d} = \{p_1, p_2, \dots, p_{20}\}$ varies a lot so that we normalize them w.r.t. the associated 2D bounding box $B^{2d} = (c_u, c_v, w, h)$ [24]. The normalized 2D

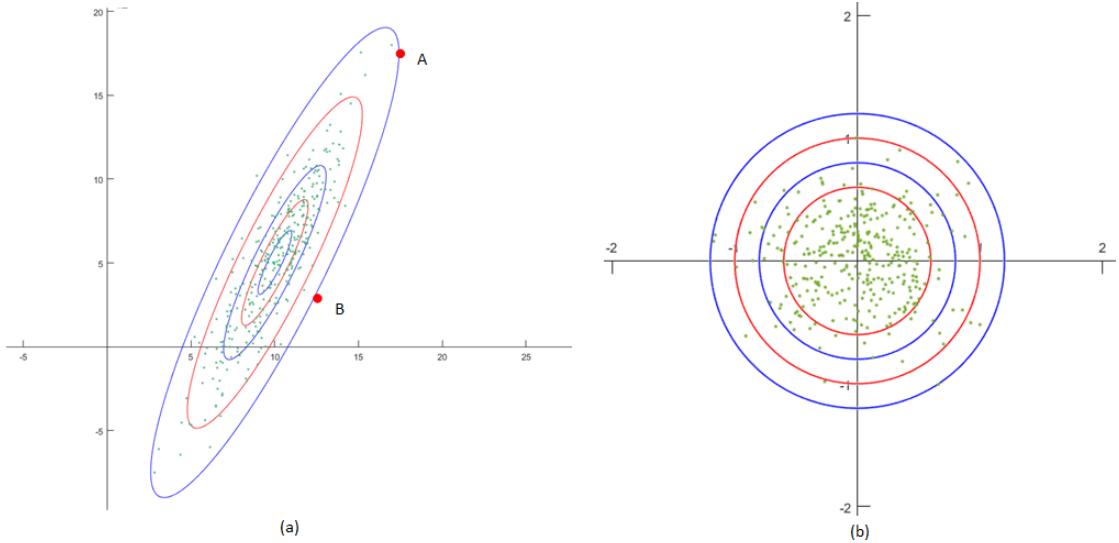


Figure 13: (a). data distribution before normalization, (b). data distribution after normalization

coordinates $\bar{C}^{2d} = \{\bar{p}_1, \bar{p}_2, \dots, \bar{p}_{20}\}$ ranges $[-1, 1]$, where

$$\bar{p}_i = \left(\frac{u_i - c_u}{w}, \frac{v_i - c_v}{h} \right) \quad (14)$$

The template proximity vector T is normalized with log function element-wisely [24], resulting in \bar{T} so that the range falls into $[-1, 1]$, too.

Moreover, as Sergey *et al.* [57] states the variation of each layer’s input distribution during the training makes the training complicated and hard to converge quickly so that we apply Batch Normalization (BN) to alleviate the internal covariate shift phenomenon. BN normalizes the summed activations of each layer to a distribution of zero mean and unit variance. The mean and variance of each activation is computed on each mini-batch. By doing so, much larger learning rate can be used for training and no special attention has to pay on parameter initialization.

4.2.2.4 Regularization

Regularization techniques are used to address the overfitting problem in order to make an algorithm not only perform well on the training data but also on the test data, the previously unseen data [48]. Overfitting results from either that the algorithm is too complicated for the data or that the sampled data is not able to represent the internal pattern. Normally, we design an algorithm to model the

data pattern exhaustively first and then apply regularization techniques to make the algorithm generalize well.

In the data aspect, we use data augmentation as a regularization. The best way to address overfitting is to train the model on more data, while the amount of data is restricted for one dataset. Therefore we generate some synthetic data. Since our network involves with classification and regression tasks, we mainly enlarge the dataset via horizontal and vertical flipping and scaling.

In the architecture aspect, Batch Normalization [57] is used as the regularization in each layer. BN enables the network to obtain the information of the training sample and the others in mini-batch simultaneously so that it does not generate deterministic dependency on this training sample. Therefore BN can replace Dropout [100] to be the regularization for our convolutional network.

Besides, Multitask Learning is another technique we used to improve generalization performance. In our network, the learning representation is shared across all tasks so that the parameters shared are constrained to bias towards one specific task, *i.e.*, only the representation that is useful for more than one can be kept [48]. Therefore statistical strength of the parameters is highly enhanced [13].

During the training, early stopping is applied to prevent the model being trained too complex to fit the test data. As Figure 14 shows, we halt the training when the validation error stops decreasing. The stopping point is when the learned model can represent the pattern of validation data most. It regards the number of training epochs as a hyperparameter and can effectively tune it to be optimal [17], because early stopping can restrict the global learning capacity of a large network to fit simpler dataset without affecting the backpropagation to control the learning capacity locally [23].

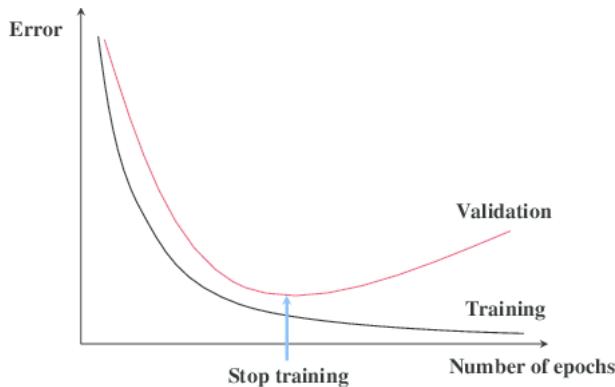


Figure 14: Early stopping

4.2.2.5 Parameter Initialization

Transfer learning is proposed to transfer the representation learned from one task to another related task [80]. As is known that CNNs learn more abstract features in a deeper layer. It is not surprising to find that, for visual tasks, shallow layers learn low-level features, *e.g.* edges, corners, changes in lighting, *etc.*, which are shared across datasets and tasks. Moreover, Yosinski *et al.* validate [116] that initializing a network with transferred features can improve the generalization capability hugely and Yoshua *et al.* [18] confirm that this initialization put the start point closer to a local minimum than random initializations, resulting in accelerating convergence.

Therefore we initialize our feature extractor with parameters learned on ImageNet dataset [90] and the three residual blocks with parameters trained on the KITTI dataset. The fully-connected layer is initialized from a zero-mean Gaussian distribution with standard deviation 0.01 and all the output layers are initialized with zeros.

4.2.2.6 Learning Algorithms / Optimization

Optimization is a task of finding the value x in order to minimize or maximize some objective function $f(x)$. One most powerful optimization technique category in deep learning is gradient-based.

As is known, the derivative $f'(x) = \frac{dy}{dx}$ specifies how to make a small change ϵ of the input x to get the corresponding change in the function:

$$f(x + \epsilon) \approx f(x) + \epsilon f'(x). \quad (15)$$

Thus, the derivative tells the direction to minimize a function, i.e. it can point out how to change x to make a small update. One popular algorithm, gradient descent [11], makes use of the derivatives and update the input x directly as:

$$x = x - \alpha \nabla_x f(x) \quad (16)$$

where α is the learning rate, a positive real value to decide the size of an update step. The examples in Figure 15 (a) show how the update works. A deep network often consists of many layers so that back-propagation algorithm [89] is used to compute the gradient for each parameter in different layers of the objective function based on the chain rule of calculus.

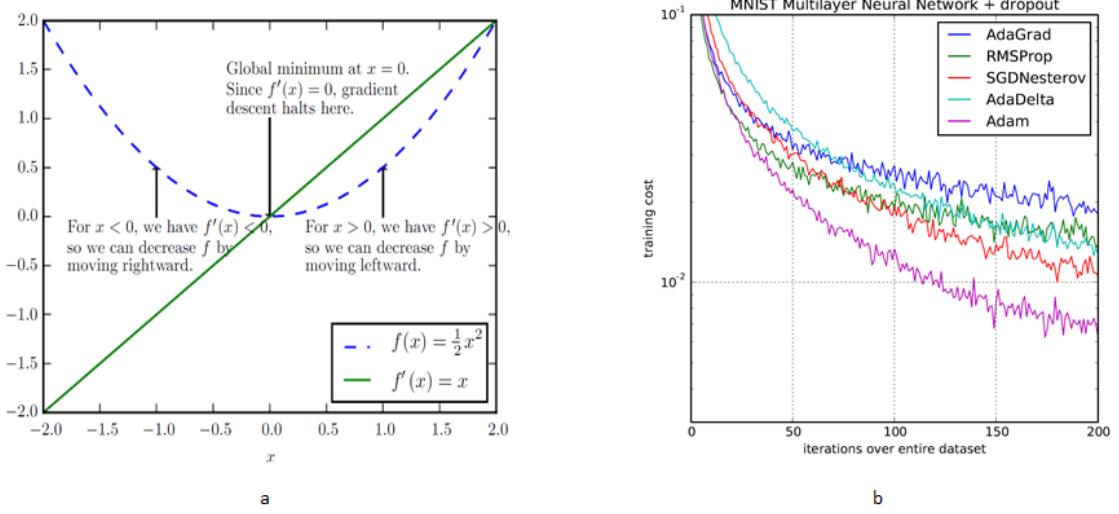


Figure 15: (a). examples to show how gradient decent makes use of derivatives to reach a minimum [48], (b). performance of optimization algorithms in the same setting [63]

In deep learning, the input to the objective function is often multidimensional so that it probably has many local minima and saddle points, which renders huge difficulties to optimization. Therefore we usually take in a compromise scenario where the value x makes f really low but not necessarily globally minimal [48].

The optimization algorithm we used is Adam [63] which is robust and efficient in memory and computation for the optimization of stochastic objectives with high-dimensional parameters spaces. It makes use of both the gradient and its momentum to update parameters as:

```

while  $x_t$  not converged do
     $t = t + 1;$ 
     $g_t = \nabla_x f(x_{t-1});$ 
     $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t;$ 
     $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2;$ 
     $\hat{m}_t = \frac{m_t}{1 - \beta_1^t};$ 
     $\hat{v}_t = \frac{v_t}{1 - \beta_2^t};$ 
     $x_t = x_{t-1} - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}};$ 
end

```

where t denotes the current iteration, α is the learning rate, β_1 and β_2 are two exponential decay rates, and ϵ is a small scalar to prevent zero-division.

From the Figure 15 (b), we can see that Adam can make the learning task converge relatively faster and has lower training error so that we follow this guide to apply Adam in our optimization of CVT Network.

4.3 Inference

As shown in Figure 16, the inference phase consists of two main steps: template matching to obtain the 3D coordinates of characteristic points and 3D dimensions for the objective vehicle and 2D-3D matching to recover the 3D vehicle location and rotation.

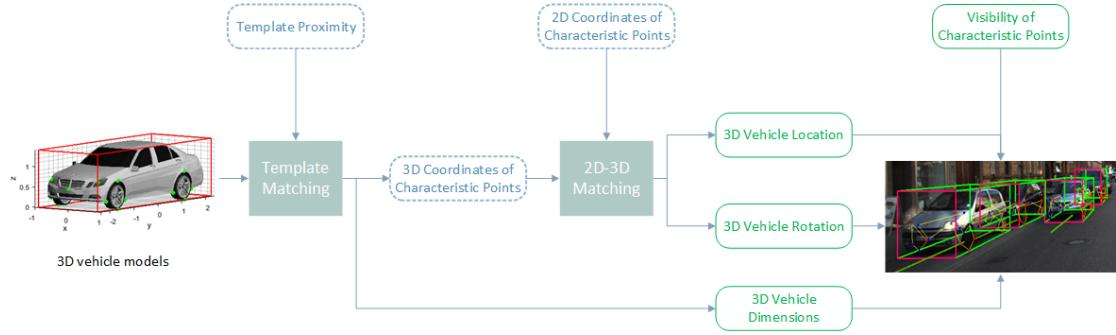


Figure 16: The Architecture of our Approach

4.3.1 Template Matching

Template matching is based on the 3D template dataset and the learned template proximity for the objective vehicle in image. As defined in section 4.1.3.3, template proximity vector, $T = \{(r_h, r_w, r_l)_k\}_{k \in \{1, \dots, K\}}$, measures the dimension similarity between the target vehicle and 3D vehicle models. The best matching model is the one whose dimensions (h, w, l) has the least distance to the target vehicle's dimension $(\bar{h}, \bar{w}, \bar{l})$, *i.e.*, the corresponding scaling ratios, (r_h, r_w, r_l) , is closest to $(1, 1, 1)$.

Let's denote the best matching template for the target vehicle m as t_j and the dimension ratio between the target vehicle and the best matching model as $r_j = (r_h, r_w, r_l)$. Then its corresponding 3D sketch is $S_j^{3d} = (p_1, p_2, \dots, p_{20})$. To get the target vehicle's dimension D_m , we apply the scaling ratios, (r_h, r_w, r_l) , to t_i as

$$D_m = t_j \cdot r_j = (h_j, w_j, l_j) \cdot (r_h, r_w, r_l) = (h^m, w^m, l^m) \quad (17)$$

In the same way, we can get the 3D coordinates of the interest points of the objective vehicle, C_m^{3d} , as

$$\begin{aligned} C_m^{3d} &= S_j^{3d} \cdot r_j \\ &= \{(x_i, y_i, z_i)\}_{i \in \{1, \dots, 20\}} \cdot (r_h, r_w, r_l) \\ &= \{(x_i^m, y_i^m, z_i^m)\}_{i \in \{1, \dots, 20\}} \end{aligned} \quad (18)$$

4.3.2 2D-3D Matching

Based on the projection mechanism described in section 3.3.2, the 2D coordinate of one point in image coordinate system can be generated via projecting the 3D point in the world coordinate system with the perspective projection matrix to the image plane. This process is described in Figure 7. Now the CVT network predicts the 2D coordinates, C_m^{2d} of the 20 interest points of the target vehicle m and the template matching process provides the corresponding 3D coordinates, C_m^{3d} . Then the 3D-2D projection equation in homogeneous coordinates becomes:

$$C_m^{2d} = K_{3 \times 4} \begin{pmatrix} R_{3 \times 3} & t_{3 \times 1} \\ 0_{1 \times 3} & 1 \end{pmatrix} C_m^{3d} \quad (19)$$

[add the optimization equation for 2D-3D matching here] where $K_{3 \times 4}$ is the given camera calibration matrix, $R_{3 \times 3}$ and $t_{3 \times 1}$ are the rotation and translation to be computed. It is easy to compute the rotation and translation of the target vehicle in the camera coordinate system with the standard 2D-3D matching algorithm, EPnP [66]. The KITTI data simplifies the rotation to one dimension, r_y , so that we follow this convention. Translation is the 3D coordinate of the origin of the vehicle so that it can represent the location of the vehicle in camera coordinate system.

TBD

4.3.3 Visualization

To visualize the 3D bounding box, we use the vehicle's dimensions to determine the eight corners of the cuboid and then project this cuboid into the image with the projection matrix which composes of the recovered rotation and translation. And in order to show the rotation clearly, project the line on the ground to show the direction of the vehicle. This visualization mechanism is shown in Figure 17 (a) where the front face of vehicle is coloured with magenta. Figure 17 also shows three typical examples in side view, back view, and front view.

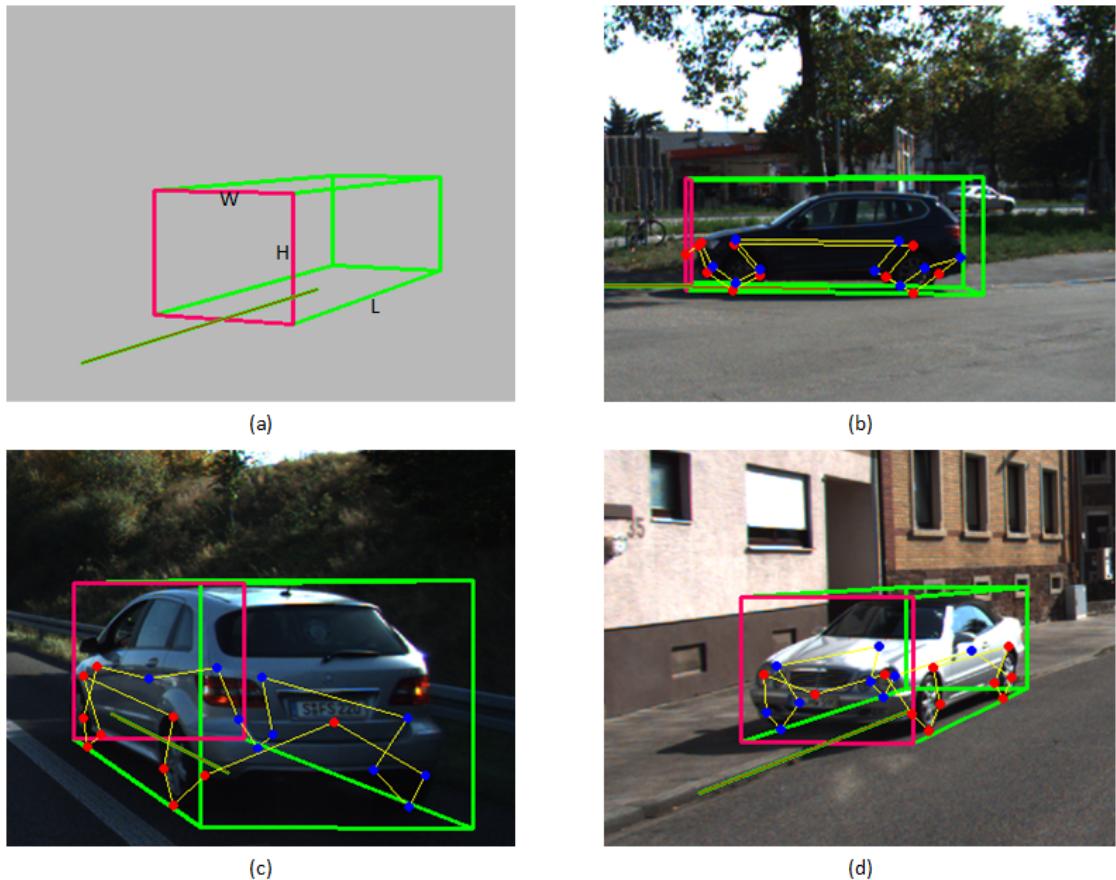


Figure 17: Visualization of 3D bounding box and visibility. (a). visualization mechanism, (b). side view of visualization, (c). back view of visualization, (d). front view of visualization

5 Experiment and Evaluation

In this section, we first describe the experiment setup and evaluation metrics in Section 5.1 and 5.2. Second we discuss the most crucial design choices for our approach in Section 5.3. Finally, we evaluate our approach and its variations , and compare our approach with the start-of-the-art methods on monocular 3D pose estimation of vehicles in Section 5.4.

5.1 Experiment Setup

We have presented the approach in Section 4.2 and 4.3 which is served as a baseline model of our approach. Some other design choices are made based on it. The benchmark network use ResNet50 [52] as feature extractor, initialized with pre-trained weights trained on ImageNet [90], and is trained with Adam [63]. The network is implemented on Keras [30] using TensorFlow [10] as backend. Keras supports running both on GPU and CPU.

We evaluate our approach and its variations on the dataset created based on KITTI 3D object detection benchmark [43], described in Section 4.1. Because the KITTI only releases the ground truth for 7481 training images, we split them into train and validation set for training and validation respectively. We follow difficulty division policy of KITTI and extend to more detailed levels. We use 54 vehicle CAD models [40] for semi-automatic labelling and template matching. Each model is encoded with 20 points for its corresponding 3D sketch.

We evaluate six tasks: 3D vehicle detection, orientation estimation, 3D localization, 3D dimension estimation, 2D part localization, and 2D part visibility prediction. 3D vehicle detection is the ultimate task, representing the localization, orientation, and dimension of 3D bounding box, which is therefore used to evaluated all the design choices.

5.2 Evaluation Metrics

We use intersection over union (IoU) to measure the performance of 3D vehicle detection. IoU used to measure the similarity of two 2D bounding boxes in various 2D object detection challenges, *e.g.* Pascal VOC [38] and ILSVRC [90]. Recently it is extended to measure 3D object detection with the formula:

$$IoU(b_1, b_2) = \frac{V(b_{pre} \cap b_{gt})}{V(b_{pre} \cup b_{gt})} \quad (20)$$

where $V(\cdot)$ indicates the volume and b_{pre} denotes the predicted 3D bounding box while b_{gt} is the ground truth. KITTI 3D object detection benchmark considers that a 3D object detection is correct, if $IoU \geq 0.7$ [43]. If multiple bounding boxes are predicted for one vehicle, they are considered as false predictions.

For 3D orientation estimation, we use the measure, Orientation Score (OS), defined in [74]. It is the mean error across all estimations in the validation set, written as:

$$OS = \frac{1}{N} \sum_{i=1}^N \frac{(1 + \cos(\Delta\theta_i))}{2} \quad (21)$$

where N denotes the number of examples in the validation set and $\Delta\theta_i$ represents the difference between the predicted orientation r_y and the ground truth for example i .

For other tasks, We follow the metrics set by [24]. A 3D localization is considered correct if its distance to the ground truth is less than a threshold. Two thresholds, 1 meter and 2 meters, are chosen. 2D part localization are measured the same way and the threshold is 20 pixels. 3D dimension estimation is correct if the predicted dimensions (h, w, l) satisfies the equation

$$\left| \frac{h - h_{gt}}{h_{gt}} \right| < 0.2 \quad \& \quad \left| \frac{w - w_{gt}}{w_{gt}} \right| < 0.2 \quad \& \quad \left| \frac{l - l_{gt}}{l_{gt}} \right| < 0.2 \quad (22)$$

where (h_{gt}, w_{gt}, l_{gt}) is the ground truth. 2D part visibility prediction is a pure classification problem so that the measure is the accuracy over 4 classes.

Task	Metric
3D vehicle detection	$IoU \geq 0.5$
3D localization	$\ \bar{\mathbf{t}}_{pre} - \bar{\mathbf{t}}_{gt}\ < 1/2$ meters
3D orientation estimation	OS
3D dimension estimation	$\left \frac{d - d_{gt}}{d_{gt}} \right _{d=\{h,w,l\}} < 20\%$
2D part localization	$\ \bar{\mathbf{p}}_{pre} - \bar{\mathbf{p}}_{gt}\ < 20$ pixels
2D part visibility	$V_{pre} = V_{gt}$

Table 4: Metrics for six tasks

5.3 Design Choices

5.3.1 Model Selection Strategy

We have mentioned in Section 4.1 that we match the vehicles with models based on dimension similarity. Now we give the reasons why we choose this strategy. The first reason is that the distribution of models' dimension covers the KITTI vehicles'. Besides, this strategy achieves the best performance among possible strategies as shown in Table 5. The CAD model dataset provides two sets of dimensions: real-world dimensions and uniform dimensions where all the width is 1.8 meters. Dims strategy means that the best model is the one whose dimension vector has the least distance to the vehicle's. Dim Ratios strategy is similar but bases on the ratio vector ($height/width, length/width$) with the uniform dimensions given by the CAD dataset. The rest two are based on Dims and Dim Ratios, *i.e.*, using one strategy to select 10 models first and then applying the other to select the best one.

We use two metrics to evaluate these strategies: type accuracy and mean deviation. The types are defined as in Figure 5. Vehicle types can present the location of characteristic points of vehicles well because the vehicles' 3D sketches varies a lot among different types but keeps similar with the same type and the dimension variation can be solved via scaling. Therefore we label 1000 KITTI vehicles and test the type matching accuracy among these strategies. The second metric, mean deviation, is the Euclidean distance between the dimension vectors of the selected model and ground truth. Even though the performance of Dims is not that good, it is the best one we can use based on the information provided. Figure 18 shows one labelled example where Dims has the most accurate points labelled and the others all label the points around the back wheel wrongly. [is it ok to put these bad results here] TBD

	Dim Ratios	Dims	Dim Ratios & Dims	Dims & Dim Ratios
Type Accuracy	0.256	0.463	0.137	0.160
Mean Deviation (m)	0.631	0.384	0.865	0.893

Table 5: Performance of four model selection strategies

5.3.2 Difficulty Levels For Dataset Division

KITTI set three levels for 3D object detection based on three factors: the size of the 2D bounding boxes, degree of occlusion and truncation. In order to compare



Figure 18: One labelled example of four model selection strategies. (a). Dim Ratios. (b). Dims. (c). Dim Ratios & Dims. (d). Dims & Dim Ratios.

the influence brought by these three factors, we extend the division into nine levels, as shown in Table 6. And Level 4 is used to evaluate other design choices because it is neither too complicate nor too simple. [more detailed level to check the difficulty of each property] TBD

Difficulty level	Min. bounding box height	Max. Occlusion	Max. Occlusion
Level 1	40 Px	0	0
Level 2 (Easy)	40 Px	0	15%
Level 3	40 Px	0, 1	15%
Level 4	40 Px	0, 1	30%
Level 5	32 Px	0, 1	30%
Level 6 (Moderate)	25 Px	0, 1	30%
Level 7	25 Px	0, 1, 2	30%
Level 8 (Hard)	25 Px	0, 1, 2	50%
Level 9	25 Px	0, 1, 2, 3	70%

Table 6: Levels of difficulty. Level 2, 6, and 8 corresponds to the Easy, Moderate, and Hard defined by KITTI [43].

We train the network for these nine levels to their best with the identical foundation

	3D bounding box estimation	3D localization (Loc. < 1m)	3D orientation estimation	3D dimension estimation	2D part localization	2D part visibility
Level 1	33.54	65.97	99.77	99.63	99.78	98.03
Level 2	34.00	67.77	99.76	99.76	99.32	97.77
Level 3	32.49	66.81	98.96	99.79	98.91	95.50
Level 4	30.92	65.60	99.01	99.73	98.72	94.65
Level 5	29.41	59.91	98.67	99.73	98.65	94.53
Level 6	26.68	57.16	98.55	99.60	98.90	94.92
Level 7	24.10	52.70	97.83	99.54	98.48	94.24
Level 8	22.86	51.35	97.30	99.48	97.52	93.92
Level 9	20.05	47.75	96.78	99.21	96.05	92.59

Table 7: Performance of six tasks at different difficult levels.

and the results of six tasks are shown in Table 7. [further analysis ref. to](#) TBD papers

5.3.3 Cost Function

We have described the loss function we used in Section 4.2.2.2 and here we give reasons for our choices.

The categorical cross entropy loss is the standard choice for multi-class classification but for regression, there are plenty of choices. We select out and do experiment on the following loss functions: mean squared error [107], mean absolute error [110], Huber loss [56], Charbonnier loss [25], a robust smooth L_1 loss function defined in [44], and our modified smooth L_1 loss function defined in Section 4.2.2.2.

We experiment with these loss functions with the same setup on Level 4 and evaluate their performance on six tasks with the validation data. The results are shown in Table 8. Our modified loss for regression tasks (the first five) outperforms all the others and only has a slight impact on the classification task. The main reason for the modification is that we normalize the labels into the range $[1, 1]$ and most of them falls in $[-0.5, 0.5]$.

Loss function	3D bounding box estimation	3D localization (Loc. < 1m)	3D orientation estimation	3D dimension estimation	2D part localization	2D part visibility
Charbonnier loss $\epsilon=0.01$	28.94	58.57	98.30	99.59	98.09	93.95
Charbonnier loss $\epsilon=0.0001$	27.3	59.04	97.89	99.73	97.7	92.27
Huber loss $\delta=0.5$	25.87	58.77	98.42	99.59	98.06	94.54
mean absolute error	4.3	17.54	90.02	92.42	86.59	93.24
mean squared error	7.51	25.94	94.02	97.61	88.81	93.48
robust smooth L_1	26.21	61.09	98.47	99.66	98.11	94.66
modified smooth L_1	27.39	62.59	98.67	99.73	98.3	94.45

Table 8: Performance of six tasks with different loss functions.

Based on the experiment above, we tune the loss weights: λ_{coord} , λ_{temp} and λ_{visib} , defined in Eq 13. As mentioned in Section 4.2.2.3, we have normalized the magnitude of all labels into the same range $[-1, 1]$. This label normalization alleviates the difficulties of loss weights tuning because we can focus on the importance of tasks without paying special attention to their quantities.

We apply grid search to find the optimal loss weight combination which results in best performance on the six tasks. Grid search is a common but sometimes expensive practice [48] but the time to train our network is acceptable because we use pre-trained models to initialize our network, . The λ_{visib} is assigned 1, λ_{coord} and λ_{temp} are chosen from the set $\{1, 3, 10, 30\}$ because 2D coordinates and template proximity is much more important than visibility for 3D object detection. The results are shown in Table 9. The ($\lambda_{coord} = 10$, $\lambda_{temp} = 1$) combination is optimal among all choices *i.e.*, it performs better on four tasks and only slightly worse on the other two tasks than other weights combinations. This is because the 2D coordinates are used to perform 2D-3D matching and a small deviation of one points in image can lead to a large difference in the world coordinate system. Besides, the template proximity can't diverge too much due to the matching strategy where the a big deviation only happens when all the predicted dimension ratios vectors are ridiculously wrong.

$(\lambda_{coord}, \lambda_{temp})$	3D bounding box estimation	3D localization (Loc. < 1m)	3D orientation estimation	3D dimension estimation	2D part localization	2D part visibility
(1, 1)	29.35	63.89	98.53	99.66	98.39	94.98
(1, 3)	25.19	56.93	97.87	99.59	97.7	94.18
(1, 10)	24.71	53.65	98.05	99.66	97.61	94.41
(1, 30)	19.59	45.39	97.68	99.59	96.46	93.92
(3, 1)	29.28	61.84	98.78	99.59	98.56	94.71
(3, 3)	29.76	63.89	98.63	99.66	98.55	94.72
(3, 10)	25.67	58.7	98.04	99.52	97.56	94.28
(3, 30)	21.98	49.69	98.32	99.59	97.18	93.99
(10, 1)	34.06	68.46	99.16	99.66	98.98	95.18
(10, 3)	30.58	63.41	98.71	99.73	98.61	94.38
(10, 10)	32.7	66.28	98.60	99.59	98.62	94.52
(10, 30)	27.03	60.27	98.09	99.73	97.79	93.96
(30, 1)	33.17	66.96	99.24	99.66	98.95	94.91
(30, 3)	31.4	66.62	98.87	99.73	98.88	94.43
(30, 10)	31.74	68.12	98.69	99.66	98.61	94.46
(30, 30)	31.19	65.32	98.66	99.52	98.36	93.43

Table 9: Performance of six tasks with different loss weights combinations.

5.3.4 Learning Rate and Weight Decay

Learning rate is one hyperparameter that determines how much the parameters of the network are updated with respect to the gradients of the cost function. It is perhaps the most crucial hyperparameter because it plays a more complicated and important role in affecting the effective capability of the model than the others. There is a typical U-shaped relationship between the learning rate and the model’s performance. If it is too small, the training is more stable but the time to train the model is much longer and sometimes, the training might get stuck on a plateau region or around a saddle point. On the other hand, if it is too large, the training can converge more quickly but is more variable because it may fail to converge or even diverge. Only a proper learning rate can achieve the best performance of the model. Therefore as Goodfellow *et al.* suggest that if you can only tune one hyperparameter, tune the learning rate [48].

Automatic learning rate selection methods alleviate the difficulty hugely, but they are too computational expensive so that we manually tune it. We first apply a “LR range test” [95], training the network for several epochs with the learning rate increasing linearly in a guessing range and then analysing the loss over these learning rates. From the test result in Figure 19, the correct learning falls in the range $[10^{-6}, 10^{-3}]$.

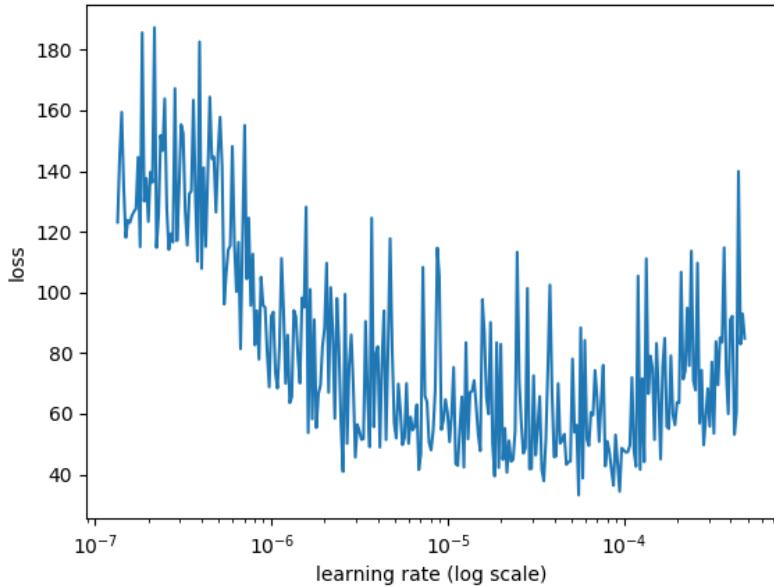


Figure 19: Results of "LR range test": loss over monotonically increasing learning rate.

Weight decay is a regularization technique which is used to prevent the weights of the model getting too big as to be overfitting. Since Keras Optimizers [30] take both learning rate and weight decay as input parameters, we apply grid search for both of them together. The results are shown in Table 10. The combination of $(10^{-5}, 10^{-6})$ for learning rate and weight decay respectively achieves the best performance except for 2D part visibility classification with 0.2% lower. Therefore we select this combination. The best results in Table 10 is worse than in Table 9 because we set all loss weights as 1 here so that they are close to the results in first row of Table 9. [wait for new results of lr5 & lr6](#)

TBD

5.3.5 Feature Extractor

Our baseline network uses ResNet50 as feature extractor. But there are several other out-of-shelf network with pre-trained weights so that we carry out experiments on them. The other base nets we choose are VGG16, VGG19 [93], InceptionV3 [102], Xception [31], DenseNet121, DenseNet169, DenseNet201[55].

As before, we train all the models with different base nets on the same setup with data Level 4. Their performance on six tasks are shown in Table 11. According to

(learning rate, weight decay)	3D bounding box estimation	3D localization (<i>Loc. < 1m</i>)	3D orientation estimation	3D dimension estimation	2D part localization	2D part visibility
($10^{-4}, 10^{-4}$)	0.491	0.564	0.977	0.997	0.981	0.948
($10^{-4}, 10^{-5}$)	0.54	0.608	0.98	0.996	0.983	0.951
($10^{-4}, 10^{-6}$)	0.499	0.566	0.978	0.997	0.984	0.951
($10^{-4}, 10^{-7}$)	0.522	0.595	0.979	0.997	0.982	0.95
($10^{-5}, 10^{-4}$)	0.519	0.596	0.979	0.994	0.975	0.932
($10^{-5}, 10^{-5}$)	0.552	0.625	0.985	0.995	0.981	0.941
($10^{-5}, 10^{-6}$)	0.539	0.631	0.986	0.997	0.986	0.949
($10^{-5}, 10^{-7}$)	0.538	0.627	0.984	0.997	0.983	0.946
($10^{-6}, 10^{-4}$)	0.382	0.442	0.962	0.99	0.951	0.892
($10^{-6}, 10^{-5}$)	0.405	0.479	0.971	0.99	0.958	0.905
($10^{-6}, 10^{-6}$)	0.38	0.46	0.969	0.988	0.952	0.897
($10^{-6}, 10^{-7}$)	0.392	0.474	0.969	0.99	0.955	0.901

Table 10: Performance of six tasks with different learning rate and weight decay combinations

their performance, only VGG19 can perform better than ResNet50 on the whole. Therefore, we includes VGG19 for our final evaluation in Section 5.4.

[is it necessary to analyse why they perform this way]

TBD

Base Net	3D bounding box estimation	3D localization (<i>Loc. < 1m</i>)	3D orientation estimation	3D dimension estimation	2D part localization	2D part visibility
ResNet50	30.85	64.98	99.03	99.73	98.77	94.74
VGG16	29.22	57.82	98.54	99.86	98.44	95.76
VGG19	31.54	64.23	99.22	99.93	98.68	96.05
InceptionV3	27.51	63.07	97.90	99.25	97.82	93.83
Xception	21.09	49.83	96.05	98.91	95.22	89.17
DenseNet121	26.28	59.73	98.15	99.52	97.47	93.26
DenseNet169	30.24	65.05	98.27	99.59	98.05	93.8
DenseNet201	27.3	61.37	98.69	99.59	98.08	93.45

Table 11: Performance of six tasks with different base net.

5.3.6 Selection of Characteristic Points

The number of characteristic points is another factor that influences the performance. Theoretically, the more points the better accuracy of 2D-3D matching because it is more robust to outliers. The outliers are eliminated via RANSAC [41]. But the labelling of these points is time-consuming and demanding. Moreover, the exact locations of points also affect the performance, too. Some points

may be hard to detect correctly with the network due to their context in the image and some even distribute differently among various shapes of vehicles.

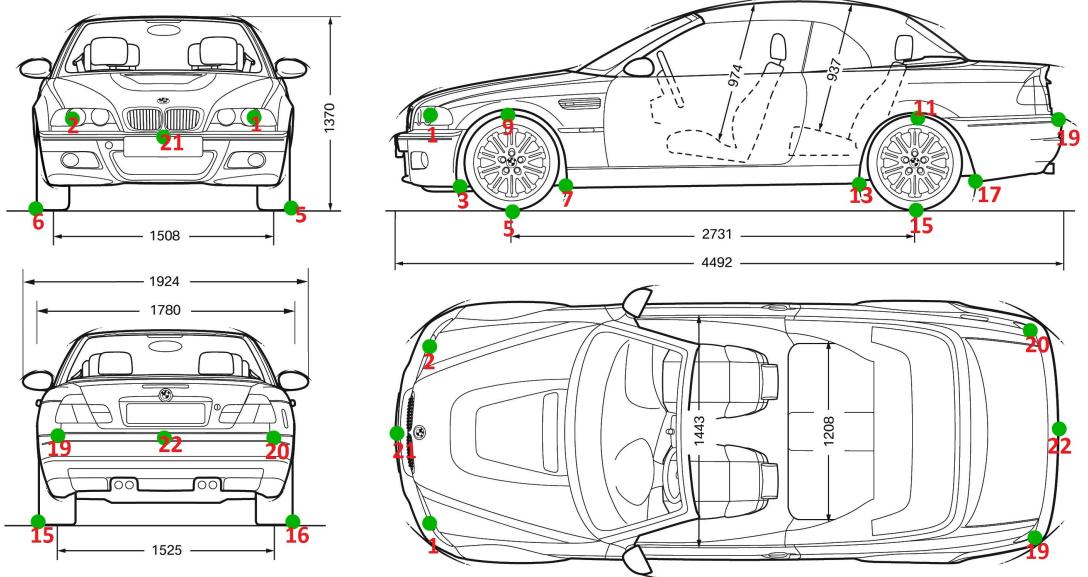


Figure 20: The schematic diagram for 22 characteristic points. point 1 and 2 denotes the center of two headlights; point 3-18 are around four wheels; point 19 and 20 denote two rear corners; and point 21 and 22 denotes the frontmost and backmost points along the ordinate axis of the vehicle.

As the Figure 20 shows, we initially label 22 points for each model to construct its corresponding 3D sketch. The points around four wheels are most important because their geometric characters are most stable and rigid, and usually, half of them are visible if not occluded by other objects. However, points 1, 2, 19, and 20 around the front and rear may vary a lot with respect to different shape of the vehicles. We still consider them because they have clear semantic meaning and we may need these points when the points around wheels are not visible, *e.g.* when the vehicle is following another vehicle on the straight lane. The geometric property of point 21 and 22 are stable but their context is not as clear as the points around the wheels.

Therefore we do some experiment on the selection of these points. We select three cases: 1). 16 points: points 3-18; 2). 20 points: points 1-18 and points 21-22; 3). 22 points: point 1-22. We train these three case on data Level 2, 6, and 8 which corresponds to the easy, moderate, and hard defined by KITTI. The reason for training on three levels instead of only Level 4 as before is that the size of the 2D bounding box and degree of the occlusion and truncation impact the geometric appearance of vehicles in images and thus affect detection of these points.

The experiment results are shown in Table 12. According to the results, 22-point case works best when the input patch is relatively bigger, *i.e.*, the minimal height for a 2D bounding box is 40 pixels. This is because when the vehicles are large in image, it is easier to estimate the 2D coordinates and the template proximity and in this scenario, the more the points, the better the performance. However, when very small vehicles are taken into consideration, it become harder to predict the 2D coordinates correctly because these key points in a small vehicle locate very close to each other.

[it requires more results to render the reasons. therefore i train them in 9 levels TBD and then to evaluate the res.]

	3D bounding box estimation	3D localization (Loc. < 1m)	3D orientation estimation	3D dimension estimation	2D part localization	2D part visibility
Easy_P16	31.28	64.81	99.75	99.88	99.48	97.82
Easy_P20	34.00	67.77	99.76	99.76	99.32	97.77
Easy_P22	34.83	69.43	99.16	99.76	99.35	97.73
<hr/>						
Moderate_P16	25.39	56.22	99.05	99.64	98.88	94.98
Moderate_P20	26.68	57.16	98.55	99.60	98.90	94.92
Moderate_P22	28.16	57.21	98.86	99.64	98.89	95.02
<hr/>						
Hard_P16	25.04	53.39	98.53	99.55	98.28	94.15
Hard_P20	22.86	51.35	97.30	99.48	97.52	93.92
Hard_P22	23.12	50.63	97.63	99.35	98.01	94.00

Table 12: The performance of six tasks with different data difficulty level and different point selection cases.

5.4 Experiment Results

In this subsection, we present the best results of six tasks achieved by our approach and compare them with the state-of-the-art works in monocular 3D vehicle detection.

5.4.1 3D Bounding Box Estimation

3D vehicle detection refers to the 3D bounding box estimation of the target vehicles in images. A 3D bounding box is estimated correctly, if $IoU > 0.7$. Our approach is based on 2D object detection, so in order to be comparable with other methods,

we eliminate the influence of this foundation by multiplying our results with the top accuracy of 2D car detection in KITTI achieved by iDST-VC [2]. Our results and the comparison are shown in Table 13. A LiDAR-based method, VoxelNet++ [119], achieves the best performance currently. Based on the assumption above, our method currently ranks at the middle (24^{th}) on 3D vehicle detection benchmark [1]. It can achieve similar performance as VoxelNet basic [119]. It outperforms all traceable image-based methods and some of LiDAR-based methods. For more comparison, check [1]. Besides, our method performs 3D detection for all cars and vans on KITTI dataset while the results from other methods are for cars only. Therefore, theoretically, the performance of our approach on cars is at least not worse than on both cars and vans.

[evaluation, mine x 2D acc. because hard to get the numerical data from graph TBD
with lowest recall]

Method	Type	Easy	Moderate	Hard
3D-SSMFCNN [78]	Mono	2.39	2.28	1.52
A3DODWTDA [50]	Mono	6.76	6.45	4.87
DoBEM [117]	LiDAR	7.42	6.95	13.45
LMNetV2 [73]	LiDAR	14.75	15.24	12.85
VoxelNet basic [119]	LiDAR	29.70	24.35	23.52
VoxelNet++ [119]	LiDAR	83.13	73.66	66.20
Ours (ResNet50)	Mono	30.90	24.16	18.53
Ours (VGG19)	Mono	30.36	23.92	20.11

Table 13: Comparison of the 3D bounding box estimation on official KITTI dataset for cars (ours is for cars and vans)

5.4.2 3D Localization

3D localization refers to the estimation of the location of the vehicle’s center. We follow the evaluation metric used by [74, 24] that a 3D localization is correct if the distance between the true and predicted center is smaller than a threshold and 1 meter and 2 meters are used as thresholds. And we borrow the idea from [74] to calculate the 3D Localization accuracy, which is the ratio between Average Localization Precision (ALP) and Average Precision (AP), for other methods. We present our results and the comparison in Table 14. Clearly, our method outperforms all other monocular methods, while can’t match the perform of stereo 3DOP [27].

[3DL = ALP / AP] [3DOP is very high, include it ?]

TBD

TBD

Method	Type	Loc. <1m			Loc. <2m		
		Easy	Moderate	Hard	Easy	Moderate	Hard
3DOP [27]	Stereo	83.69	65.96	65.71	98.57	88.08	91.56
DPM [39]	Mono	34.33	29.02	29.18	56.48	46.69	46.17
3DVP [112]	Mono	52.15	45.24	42.40	76.10	68.00	64.84
SubCNN [113]	Mono	43.26	34.86	32.75	77.66	63.12	59.33
Ours (ResNet50)	Mono	67.87	59.96	51.38	87.56	79.83	73.11
Ours (VGG19)	Mono	68.25	56.58	54.53	88.39	78.63	77.49

Table 14: Comparison of the 3D Localization accuracy on official KITTI dataset for cars (ours is for cars and vans)

5.4.3 3D Orientation Estimation

3D Orientation Estimation refers to the estimation of rotation of vehicle along the Y-axis (yaw axis), r_y , in the camera coordinate systems. The evaluation metric is Orientation Score (OS), which is the ratio between Average Orientation Estimation (AOS) and Average Precision (AP) [74]. Table 15 shows the results of our methods and other state-of-the-art ones. Although our method can't achieve the better performance than , the performance than SubCNN [113], Deep3DBox [74], and DeepMANTA [24] it performs better than 3DOP [27] and Mono3D [26]. Therefore, our method achieves the start-of-the-art level performance for 3D orientation estimation. Besides, our approach consider both cars and vans from KITTI dataset, while the others are only evaluated on cars.

[OS = AOS / AP]

TBD

Method	Type	Easy	Moderate	Hard
3DOP [27]	Stereo	98.28	97.13	96.73
Mono3D [26]	Mono	98.57	97.69	97.31
SubCNN [113]	Mono	99.84	99.52	99.25
Deep3DBox [74]	Mono	99.91	99.67	99.46
DeepMANTA (GoogLeNet) [24]	Mono	99.85	99.75	99.55
DeepMANTA (VGG16) [24]	Mono	99.88	99.69	99.49
Ours (ResNet50)	Mono	99.88	98.78	97.39
Ours (VGG19)	Mono	99.94	99.09	98.43

Table 15: Comparison of the Orientation Score (OS) on official KITTI dataset for cars (ours is for cars and vans).

5.4.4 3D dimension estimation, 2D part localization, and 2D part visibility

Our approach is developed based on DeepMANTA [24], so only DeepMANTA shares these three tasks with us. We present the results of DeepMANTA and our approach in Table 16. And ours outperforms DeepMANTA on all these three tasks at all difficulty levels.

Method	Type	3D dimension estimation			2D part localization			2D part visibility		
		Easy	Moderate	Hard	Easy	Moderate	Hard	Easy	Moderate	Hard
DeepMANTA [24]	Mono	97.54	90.79	82.64	92.48	85.08	76.9	94.04	86.62	78.72
Ours (ResNet50)	Mono	99.76	99.6	99.48	99.32	98.9	97.52	97.77	94.92	93.92
Ours (VGG19)	Mono	100	99.78	99.77	99.96	99.21	98.4	98.36	96.06	95.52

Table 16: Comparison of 3D dimension estimation, 2D part localization, and 2D part visibility on official KITTI dataset for cars (ours is for cars and vans).

[the methods used to compare on different tasks are inconsistent] [better to show some images here to show the performance]

TBD

TBD

6 Discussion

In this section, we discuss the deficiencies of our approach, analyse their causes and propose some possible improvements.

6.1 Deficiency of 2D Coordinate Labelling

6.1.1 Deficiency of Our Labelling Approach

Our approach requires accurate 2D characteristic points to perform 2D-3D matching to calculate the location and orientation of the vehicles. However precisely labelling these points is demanding or even impossible for some vehicles with our semi-automatic labelling tool. The deficiency of labelling results from our labelling approach and the KITTI dataset.

Due to the variety of dimensions and shapes of vehicles, it is very hard to find a model that exactly matches the target vehicle with a small set of models. But if we expand the model set to cover all vehicles, the required quantity is too huge to satisfy. Therefore, we can only use a small set of models which means that the approximation and scaling are inevitable. And this will introduce errors to the labelling of 2D coordinates. Figure 21 (a) and (b) show two bad labelled examples. The imprecision of (a) results from the approximation of model matching, *i.e.*, most points are labelled correctly but point 1-6 are not. This is because the model can't perfectly fit the vehicle. The wrongly labelled points around wheels roots in the scaling the model to fit the vehicle so that, for example, the distance between point 4 and 8 are larger than the diameter of the wheel. The other pairs of points around the wheel have similar situation so that most points are not labelled correctly for this vehicle.

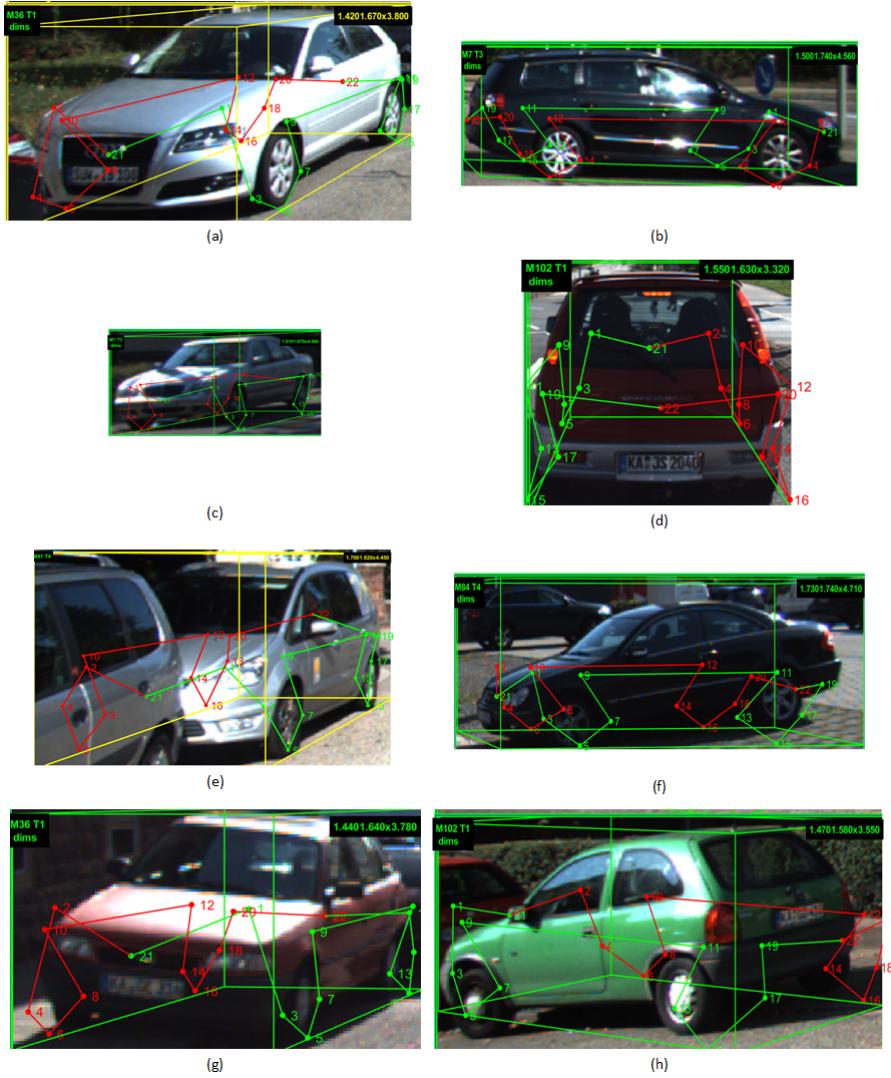


Figure 21: Labelling examples for 2D points with original size. (a). Approximation results in inaccurate label for point 1-6. (b). Scaling change the relative distance between points which lead to bad labelling for points around the wheels. (c). One of the very small vehicles in image is able to be labelled with our approach. (d). The example that the side faces are not visible is able to be labelled with our approach. (e). The occluded vehicle can also be labelled with our approach. (f). The 3D bounding box ground truth of this vehicle is not correctly labelled in KITTI dataset. (g). KITTI doesn't consider the rotation along the roll axis, which results in bad labelling. (h). KITTI doesn't consider the rotation along the pitch axis, which results in totally wrong labelling.

6.1.2 Deficiency of KITTI Dataset Labelling

Our labelling approach requires very accurate label of the dataset because 2D points are very sensitive to the exact location of the target vehicle. But some of the KITTI images are not labelled so well that our approach makes mistakes. There are three main deficiencies of the KITTI labelling as shown in Figure 21 (f), (g), and (h). (f) shows the 3D bounding box ground truth of this vehicle is not correctly labelled. The 3D bounding box doesn't encompass the vehicle, especially the rear part, which leads to wrong labelling. (g) and (h) shows that our approach wrong labels the 2D points due to the assumption of KITTI that the vehicle only rotate along the yaw axis rather than the roll and pitch axis.

6.1.3 Possible Improvements

Despite the drawbacks of our labelling approach, it has still has some advantages. Our approach can label very small vehicles in the image as shown in Figure 21 (c). Besides, it can label the self-occluded points as in Figure 21 (d) and the occluded and truncated points in Figure 21 (e).

Therefore, if labelling labour is allowed, we can manually label the visible points for big vehicles and use the projection property and the geometry constraints to automatically generate the corresponding points. Therefore we have more accurate points for big vehicles. And for small and highly occluded or truncated vehicles, it is suitable to use our approach to label them. Because manually labelling these special vehicles is almost impossible and the error of our labelling approach can be decreased due to the quantities.

6.2 Deficiency of Visibility Labelling

6.2.1 Deficiency of Our Labelling Approach

As described in Section 4.1.3.2, we label the visibility property of characteristic points with the help of the rotation r_y , other labelled objects, and the size of images. But there are two main shortcomings. The first one is that the 2D bounding box of a vehicle is larger in image than its real shape which makes our approach mistake the visible or self-occlude points as occluded. Figure 22 (a). shows an example of labelling the self-occluded points as occluded, which leads to wrong labels. Another shortcoming is that our approach can't take the unlabelled objects into consideration so that it automatically ignores than during

labelling. Figure 22 (b) and (c) show two examples for this case. Our approach can't consider the traffic sign in (b) and the building in (c) and thus, it fails to label these occluded points correctly.

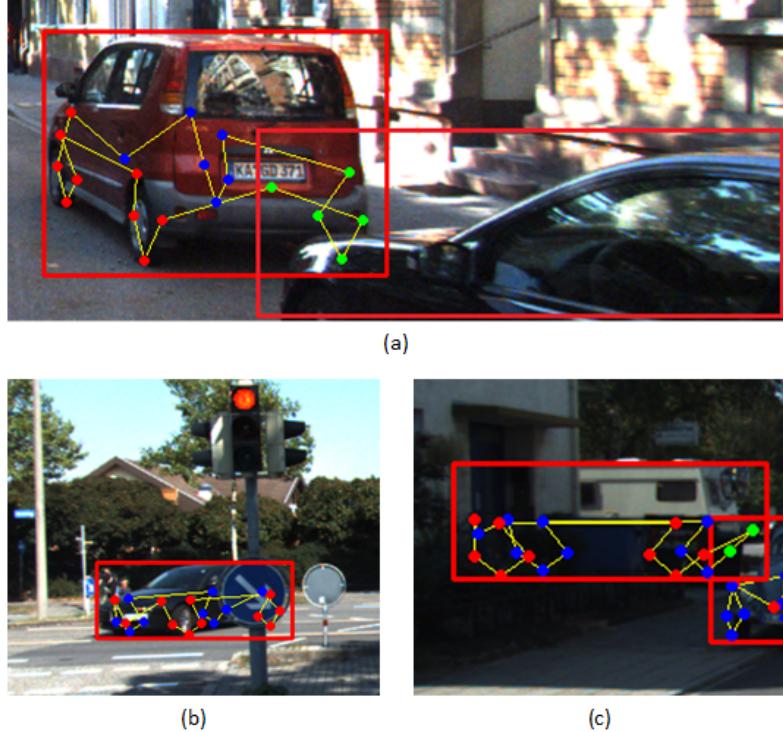


Figure 22: Labelling examples for visibility with original size. (a). The 2D bounding box is larger than the vehicle's real size, which results in labelling the self-occluded points as occluded. (b). KITTI doesn't label the traffic sign so that our approach can't label the occluded points correctly. (c). KITTI doesn't label the building so that our approach can't label the occluded points correctly.

6.2.2 Possible Improvements

KITTI offers the LiDAR data for all images, which can provide the depth information of each pixel point in image or at least of a very small region. And based on the ground truth of location, orientation, and dimension of vehicles given by KITTI, we can calculate the depth of points distributed on the surface of the vehicles. By comparing these two kinds of depth information, we can correctly figure out whether there is some object in front of the target vehicle or not and which parts of the target vehicle are occluded, and thus we can classify the visibility properties correctly. Figure 23 shows the images with LiDAR projection for exam-

ples in Figure 22. With the LiDAR depth information, we can easily classify the points around rear in (a) as self-occluded, and the occluded objects in (b) and (c) can be clearly detected and thus, we can assign the points whose depth is deeper than these objects as occluded. Therefore, for more accurate visibility labelling, it is worth making use of the LiDAR data.

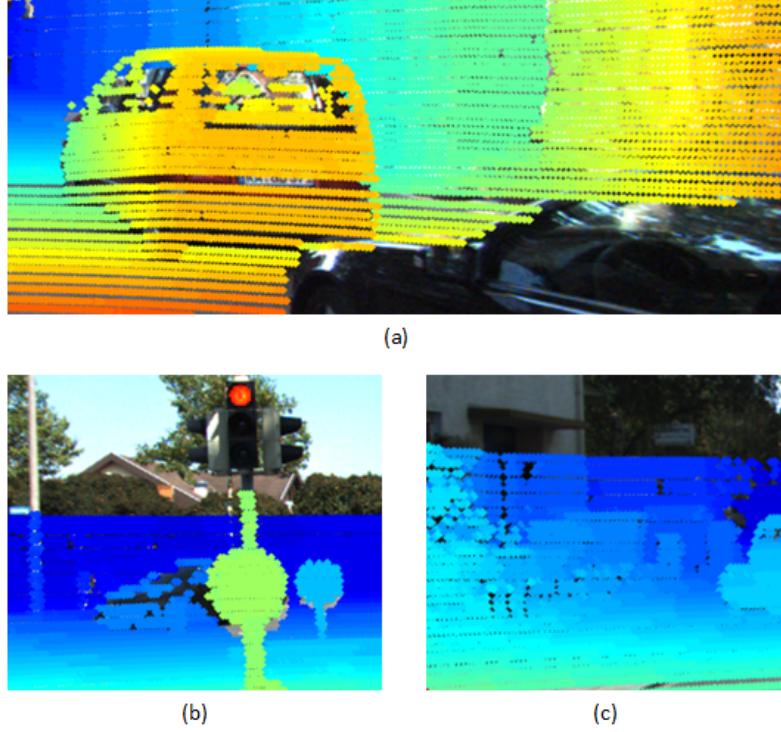


Figure 23: Images with LiDAR data projection for examples in Figure 22 (a), (b), and (c) respectively.

6.3 Deficiency of Template Matching

In our approach, there are two situations where template matching is performed to fit a target vehicle with a 3D vehicle model. The first one is during labelling, when a 3D vehicle sketch is required to be projected into the KITTI images to generate 2D coordinates of key points for a target vehicle. And the model is selected via dimension matching. Another is during inference phase, where one 3D sketch is needed to generate the corresponding 3D coordinates of the key points for the target vehicle in the world coordinate system. The time the model is selected based on the estimated template proximity vector, *i.e.*, the ratios of 3 dimensions between the target vehicle and all 3D vehicle models.

For both cases, the most fundamental drawback is that we can only match the target vehicle with an approximate model based on the Dims strategy. This would result in errors in 2D and 3D coordinates of key points and thus impacts the final 3D vehicle pose estimation.

[distribution of dims of models and kitti. labelled 23 new models but maybe no time to retain them.] TBD

6.3.1 Possible Improvements

As we mentioned before, the geometric property of 3D sketch relies heavily on the type of the vehicles. Therefore, a possible improvement for matching is to introduce type information into the matching phase. And there are some off-the-shelf frameworks that provide this kind of information [115, 96, 69, 33], *i.e.*, Afshin *et al.* provide a network that can recognize the Make and Model of vehicles [33].

In sum, the optimal matching strategy is to first collect a 3D vehicle dataset where the distribution of vehicle dimensions covers the main part of the distribution of real vehicle dimensions in a fine-grained style, and during the matching, we first estimate the type of the vehicle and then select the best model via dimension matching from the same category.

6.4 Deficiency of Input Image Format

Now our network takes fixed-size patches as inputs which are generated by extracting, zero-padding, and resizing the 2D bounding box regions and then extracts features from these wrapped patches for regression and classification tasks. The data input mode is in spirit similar to R-CNN [45], where we have to process one KITTI image multiple times to estimate all vehicles in it. The reason why we use this data input format is to make sure that the input images in one mini-batch have the same size to enjoy the efficient matrix multiplication while the 2D bounding box size is varying among different vehicles.

Even if no extra computation is wasted on the non-vehicle regions, the multiple-time input mode is not that optimal because the patch generation stage is separated with the network and it's cumbersome to test different input formats, *i.e.*, patch scales and data augmentation.

6.4.1 Possible Improvements

In the future, we plan to exploit the Fast-RCNN [44] data input mode rather than R-CNN mode. The network first takes one entire KITTI image and its 2D bounding boxes as input, generates a shared feature map for this image and then extracts a fixed-size feature vector from each feature map portion corresponding to its 2D bounding box one by one for later predictions. This functionality is realized mainly by introducing an ROI pooling layer described in [44].

Not like Fast-RCNN [44], the number of ROI, *i.e.*, bounding boxes, is dynamic for different images in our case. Therefore we have to construct a dynamic computation graph for the network during training and inference. This is why we haven't use this alternative up till now because we first implement the network in a static mode in Keras [30] with TensorFlow [10] as backend and it is hard to realize our functionality with them. Therefore, in the future, we will make use of libraries, *i.e.*, PyTorch [81], to implement this alternative.

In this way, we can process the whole image one time rather than multiple times which can improve the process speed. What's more important is that we can easily modify how large the feature map portions are used to generate a fixed-size feature vector. Based on this, we can easily evaluate whether the additional context region around the 2D bounding boxes are beneficial for 3D pose estimation and if so, how much context information can boost the performance most.

6.5 Deficiency of the Whole Approach

The first unpleasant thing of this approach is the additional labelling work. This is so time consuming that it takes more than two month to find a proper method to generate the additional labels at an acceptable accuracy level. Besides, even with this method, the accuracy of the labels is compromised.

Moreover, our approach requires an additional 3D vehicle model dataset to encode the geometry information of vehicles, which makes it hard to generalize to vehicles without a corresponding model in the model dataset. For example, our approach can't perform accurate 3D pose estimation of buses because we don't include the bus models.

According to the competition results of KITTI 3D object detection, the approach that relies on single sensor can't compare with those making use of multi-sensor data [1]. RGB images captured by cameras have high resolution and good texture knowledge but lack depth information, which cloud points collected by LiDAR have 3D information of the scene but the resolution is relative lower and lack

texture information. Therefore we think the best performance would be achieved by some approach driven by sensor fusion. If we have can start this task again, this is the first direction we will explore.

[is it necessary to analyse why our approach can't achieve very high accuracy, *i.e.*,
analyse the wrong examples] TBD

[is it necessary to point out the advantages of our approach] TBD

7 Conclusion

This thesis deals with the problem of 3D pose estimation of vehicles based on monocular images by using deep neural networks for the application of autonomous driving. We train a deep neural network to predict 2D coordinates and visibility property for characteristic points, as well as a template proximity vector, for a target vehicle. During the inference phase, we first perform template matching to reason out the target vehicle’s 3D dimensions and 3D coordinates of the key points with the help of the template proximity vector; and then 2D-3D matching is performed to recover the location and orientation of the vehicle. Combining the network and inference phases, our approach can simultaneously performs 3D vehicle detection, 3D localization, 3D orientation estimation, 3D dimension estimation, 2D part localization, and parts visibility characterization for a vehicle in a 2D bounding box patch.

Our work achieves state-of-the-art performance on six tasks. It outperforms most monocular methods recorded in KITTI 3D object detection competition on the most important task, 3D vehicle detection. Besides, our network can predict 2D coordinates and template proximity vector for highly occluded or truncated vehicles and therefore, our approach can perform template matching and 2D-3D matching to recover their 3D bounding boxes. Moreover, unlike other methods, our approach can provide more detailed information of the detected vehicle, *e.g.* 2D part location and 2D part visibility, which are useful for autonomous driving applications to gain a more fine-grained perception. Finally, the runtime for our approach is at real time level, *ca.* 0.02s per image on one GeForce GTX TITAN X GPU.

In order to maximize the effective capacity of the model, we thoroughly research and evaluate the key design alternatives and hyperparameters, including loss function, loss weight, learning rate, weight decay, base net feature extractor, model selection strategy, and number of points used in a model. Based on these experiments, the performance of our approach is boosted steadily and finally reaches the state-of-the-art level.

Finally, we present results of these six tasks and compare them with other methods. We also describe the deficiencies of our approach, analyse their causes, and propose possible solutions which are the directions of future improvements.

References

- [1] 3d object detection evaluation 2017.
- [2] The best performance of 2d car detection in kitti.
- [3] Data preprocessing.
- [4] History of autonomous cars. https://en.wikipedia.org/wiki/History_of_autonomous_cars#cite_note-1061.
- [5] No hands across america home page. http://www.cs.cmu.edu/afs/cs/usr/tjochem/www/nhaa/nhaa_home_page.html.
- [6] One-hot-encoding.
- [7] What is computer vision.
- [8] Phantom auto will tour city. *The Milwaukee Sentinel. Google News Archive*, Dec 1926.
- [9] *Taxonomy and Definitions for Terms Related to On-Road Motor Vehicle Automated Driving Systems*, Jan 2014.
- [10] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [11] Cauchy Augustin-Louis. Méthode générale pour la résolution des systèmes d'équations simultanées. *Compte rendu des séances de l'académie des sciences*, abs/1411.1792:536–538, 1847.
- [12] Dana H. Ballard, Geoffrey E. Hinton, and Terrence J. Sejnowski. Parallel visual computation. *Nature*, 306:21–26, 1983.
- [13] Jonathan Baxter. Learning internal representations. In *Proceedings of the Eighth Annual Conference on Computational Learning Theory*, COLT '95, pages 311–320, New York, NY, USA, 1995. ACM.

- [14] Tal Ben-Nun and Torsten Hoefer. Demystifying parallel and distributed deep learning: An in-depth concurrency analysis. *CoRR*, abs/1802.09941, 2018.
- [15] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, Aug 2013.
- [16] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *Trans. Neur. Netw.*, 5(2):157–166, March 1994.
- [17] Yoshua Bengio. Practical recommendations for gradient-based training of deep architectures. *CoRR*, abs/1206.5533, 2012.
- [18] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. Greedy layer-wise training of deep networks. In B. Schölkopf, J. C. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 153–160. MIT Press, 2007.
- [19] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [20] Y-Lan Boureau, Jean Ponce, and Yann LeCun. A theoretical analysis of feature pooling in visual recognition. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ICML’10, pages 111–118, USA, 2010. Omnipress.
- [21] Alberto Broggi, Pietro Cerri, Mirko Felisa, Maria Chiara Laghi, Luca Mazzei, and Pier Paolo Porta. The vislab intercontinental autonomous challenge: an extensive test for a platoon of intelligent vehicles. *International Journal of Vehicle Autonomous Systems*, 10(3):147–164, 2012.
- [22] Rich Caruana. Multitask learning. *Machine Learning*, 28(1):41–75, Jul 1997.
- [23] Rich Caruana, Steve Lawrence, and Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Proceedings of the 13th International Conference on Neural Information Processing Systems*, NIPS’00, pages 381–387, Cambridge, MA, USA, 2000. MIT Press.
- [24] Florian Chabot, Mohamed Chaouch, Jaonary Rabarisoa, Céline Teulière, and Thierry Chateau. Deep MANTA: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image. *CoRR*, abs/1703.07570, 2017.

- [25] P. Charbonnier, L. Blanc-Feraud, G. Aubert, and M. Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *Proceedings of 1st International Conference on Image Processing*, volume 2, pages 168–172 vol.2, Nov 1994.
- [26] Xiaozhi Chen, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. Monocular 3d object detection for autonomous driving. In *IEEE CVPR*, 2016.
- [27] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Andrew Berneshawi, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals for accurate object class detection. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’15, pages 424–432, Cambridge, MA, USA, 2015. MIT Press.
- [28] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals using stereo imagery for accurate object class detection. *CoRR*, abs/1608.07711, 2016.
- [29] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. *CoRR*, abs/1611.07759, 2016.
- [30] François Chollet et al. Keras. <https://keras.io>, 2015.
- [31] François Chollet. Xception: Deep learning with depthwise separable convolutions. *CoRR*, abs/1610.02357, 2016.
- [32] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-FCN: object detection via region-based fully convolutional networks. *CoRR*, abs/1605.06409, 2016.
- [33] Afshin Dehghan, Syed Zain Masood, Guang Shu, and Enrique G. Ortiz. View independent vehicle make, model and color recognition using convolutional neural network. *CoRR*, abs/1702.01721, 2017.
- [34] V. Dhiman, Q. H. Tran, J. J. Corso, and M. Chandraker. A continuous occlusion model for road scene understanding. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4331–4339, June 2016.
- [35] Xinxin Du, Marcelo H. Ang Jr., Sertac Karaman, and Daniela Rus. A general pipeline for 3d detection of vehicles. *CoRR*, abs/1803.00387, 2018.
- [36] Andreas Eitel, Jost Tobias Springenberg, Luciano Spinello, Martin A. Riedmiller, and Wolfram Burgard. Multimodal deep learning for robust RGB-D object recognition. *CoRR*, abs/1507.06821, 2015.

- [37] Martin Engelcke, Dushyant Rao, Dominic Zeng Wang, Chi Hay Tong, and Ingmar Posner. Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks. *CoRR*, abs/1609.06666, 2016.
- [38] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, January 2015.
- [39] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(9):1627–1645, Sept 2010.
- [40] Sanja Fidler, Sven Dickinson, and Raquel Urtasun. 3d object detection and viewpoint estimation with a deformable 3d cuboid model. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 611–619. Curran Associates, Inc., 2012.
- [41] Martin A. Fischler and Robert C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, June 1981.
- [42] David A. Forsyth and Jean Ponce. *Computer Vision: A Modern Approach*. Prentice Hall Professional Technical Reference, 2002.
- [43] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [44] Ross B. Girshick. Fast R-CNN. *CoRR*, abs/1504.08083, 2015.
- [45] Ross B. Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. *CoRR*, abs/1311.2524, 2013.
- [46] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterington, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.

- [47] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 315–323, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.
- [48] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [49] Saurabh Gupta, Ross B. Girshick, Pablo Arbelaez, and Jitendra Malik. Learning rich features from RGB-D images for object detection and segmentation. *CoRR*, abs/1407.5736, 2014.
- [50] Fredrik Gustafsson and Erik Linder-Norén. Automotive 3d object detection without target domain annotations. Master’s thesis, Linköping University, 2018.
- [51] Kaiming He and Jian Sun. Convolutional neural networks at constrained time cost. *CoRR*, abs/1412.1710, 2014.
- [52] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [53] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [54] Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *CoRR*, abs/1207.0580, 2012.
- [55] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks. *CoRR*, abs/1608.06993, 2016.
- [56] Peter J. Huber. Robust estimation of a location parameter. *Ann. Math. Statist.*, 35(1):73–101, 03 1964.
- [57] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR*, abs/1502.03167, 2015.
- [58] Kevin Jarrett, Koray Kavukcuoglu, Karol Gregor, and Yann LeCun. What is the best feature learning procedure in hierarchical recognition architectures? *CoRR*, abs/1606.01535, 2016.
- [59] T Jayalakshmi and Santhakumaran A. Statistical normalization and back propagation for classification. 3:89–93, 01 2011.

- [60] S. Buehler M. Iagnemma K., Singh. *The 2005 DARPAvGrand Challenge: The Great Robot Race*, volume 36. Springer, 1st. edition, 2007.
- [61] S. Buehler M. Iagnemma K., Singh. *The DARPA Urban Challenge: Autonomous Vehicles in City Traffic*, volume 56. Springer, 1st. edition, 2009.
- [62] Takeo Kanade, Chuck Thorpe, and William Whittaker. Autonomous land vehicle project at cmu. In *Proceedings of the 1986 ACM Fourteenth Annual Conference on Computer Science*, CSC '86, pages 71–80, New York, NY, USA, 1986. ACM.
- [63] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [64] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven Lake Waslander. Joint 3d proposal generation and object detection from view aggregation. *CoRR*, abs/1712.02294, 2017.
- [65] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998.
- [66] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. Epnp: An accurate o(n) solution to the pnp problem. *International Journal of Computer Vision*, 81(2):155, Jul 2008.
- [67] Bo Li. 3d fully convolutional network for vehicle detection in point cloud. *CoRR*, abs/1611.08069, 2016.
- [68] Bo Li, Tianlei Zhang, and Tian Xia. Vehicle detection from 3d lidar using fully convolutional network. *CoRR*, abs/1608.07916, 2016.
- [69] T. Y. Lin, A. RoyChowdhury, and S. Maji. Bilinear cnn models for fine-grained visual recognition. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1449–1457, Dec 2015.
- [70] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. Feature pyramid networks for object detection. *CoRR*, abs/1612.03144, 2016.
- [71] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *CoRR*, abs/1708.02002, 2017.
- [72] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: single shot multibox detector. *CoRR*, abs/1512.02325, 2015.

- [73] K. Minemura, H. Liau, A. Monrroy, and S. Kato. LMNet: Real-time Multiclass Object Detection on CPU using 3D LiDAR. *ArXiv e-prints*, May 2018.
- [74] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. *CoRR*, abs/1612.00496, 2016.
- [75] Xin C. Wang Z. Zhang N. Zheng, J. *China future challenge: Beyond the intelligent vehicle*, volume 16, pages 8–10. IEEE Intell. Transp. Syst. Soc. Newslett, 2014.
- [76] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ICML’10, pages 807–814, USA, 2010. Omnipress.
- [77] Michael Nielsen. Neural networks and deep learning, Dec. 2017.
- [78] Libor Novak. Vehicle detection and pose estimation for autonomous driving. Master’s thesis, Czech Technical University in Prague, 2017.
- [79] University of Washington. Park shuttle automated driverless vehicle. <http://faculty.washington.edu/jbs/ittrans/parkshut.htm>, 2009.
- [80] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Trans. on Knowl. and Data Eng.*, 22(10):1345–1359, October 2010.
- [81] 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 *NIPS-W*, 2017.
- [82] Cuong Cao Pham and Jae Wook Jeon. Robust object proposals re-ranking for object detection in autonomous driving using convolutional neural networks. *Signal Processing: Image Communication*, 53:110 – 122, 2017.
- [83] Dean A. Pomerleau. Alvinn: An autonomous land vehicle in a neural network. In D. S. Touretzky, editor, *Advances in Neural Information Processing Systems 1*, pages 305–313. Morgan-Kaufmann, 1989.
- [84] Charles Ruizhongtai Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J. Guibas. Frustum pointnets for 3d object detection from RGB-D data. *CoRR*, abs/1711.08488, 2017.

- [85] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CoRR*, abs/1612.00593, 2016.
- [86] Joseph Redmon and Ali Farhadi. YOLO9000: better, faster, stronger. *CoRR*, abs/1612.08242, 2016.
- [87] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *CoRR*, abs/1506.01497, 2015.
- [88] Sebastian Ruder. An overview of gradient descent optimization algorithms. *CoRR*, abs/1609.04747, 2016.
- [89] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Neurocomputing: Foundations of research. chapter Learning Representations by Back-propagating Errors, pages 696–699. MIT Press, Cambridge, MA, USA, 1988.
- [90] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014.
- [91] J. Schlosser, C. K. Chow, and Z. Kira. Fusing lidar and images for pedestrian detection using convolutional neural networks. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2198–2205, May 2016.
- [92] Jürgen Schmidhuber. Prof. schmidhuber’s highlights of robot car history. <http://people.idsia.ch/~juergen/robotcars.html>, 2009.
- [93] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [94] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [95] Leslie N. Smith. No more pesky learning rate guessing games. *CoRR*, abs/1506.01186, 2015.
- [96] J. Sochor, A. Herout, and J. Havel. Boxcars: 3d boxes as cnn input for improved fine-grained vehicle recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3006–3015, June 2016.

- [97] S. Song and M. Chandraker. Joint sfm and detection cues for monocular 3d localization in road scenes. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3734–3742, June 2015.
- [98] S. Song and J. Xiao. Deep sliding shapes for amodal 3d object detection in rgbd images. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 808–816, June 2016.
- [99] Shuran Song and Jianxiong Xiao. Sliding shapes for 3d object detection in depth images. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 634–651, Cham, 2014. Springer International Publishing.
- [100] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [101] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014.
- [102] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015.
- [103] Richard Szeliski. Computer vision algorithms and applications, 2011.
- [104] Yichuan Tang. Deep learning using support vector machines. *CoRR*, abs/1306.0239, 2013.
- [105] Dominic Zeng Wang and Ingmar Posner. Voting for voting in online point cloud object detection. In *Proceedings of Robotics: Science and Systems*, Rome, Italy, July 2015.
- [106] T. Wang, D. J. Wu, A. Coates, and A. Y. Ng. End-to-end text recognition with convolutional neural networks. In *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, pages 3304–3308, Nov 2012.
- [107] Z. Wang and A. C. Bovik. Mean squared error: Love it or leave it? a new look at signal fidelity measures. *IEEE Signal Processing Magazine*, 26(1):98–117, Jan 2009.
- [108] Zining Wang, Wei Zhan, and Masayoshi Tomizuka. Fusing bird view LIDAR point cloud and front view camera image for deep object detection. *CoRR*, abs/1711.06703, 2017.

- [109] R. G. J. Wijnhoven and P. H. N. de With. Fast training of object detection using stochastic gradient descent. In *2010 20th International Conference on Pattern Recognition*, pages 424–427, Aug 2010.
- [110] C. Willmott and K. Matsuura. Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance. *Climate Research*, 30:79–82, 2005.
- [111] Bichen Wu, Alvin Wan, Xiangyu Yue, and Kurt Keutzer. Squeezeseg: Convolutional neural nets with recurrent CRF for real-time road-object segmentation from 3d lidar point cloud. *CoRR*, abs/1710.07368, 2017.
- [112] Yu Xiang, Wongun Choi, Yuanqing Lin, and Silvio Savarese. Data-driven 3d voxel patterns for object category recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1903–1911. 2015.
- [113] Yu Xiang, Wongun Choi, Yuanqing Lin, and Silvio Savarese. Subcategory-aware convolutional neural networks for object proposals and detection. *CoRR*, abs/1604.04693, 2016.
- [114] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. *CoRR*, abs/1711.10871, 2017.
- [115] Linjie Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. A large-scale car dataset for fine-grained categorization and verification. *CoRR*, abs/1506.08959, 2015.
- [116] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *CoRR*, abs/1411.1792, 2014.
- [117] S. L. Yu, T. Westfertel, R. Hamada, K. Ohno, and S. Tadokoro. Vehicle detection and localization on bird’s eye view elevation images using convolutional neural network. In *2017 IEEE International Symposium on Safety, Security and Rescue Robotics (SSRR)*, pages 102–109, Oct 2017.
- [118] Jing Zhou, Xiaopeng Hong, Fei Su, and Guoying Zhao. Recurrent convolutional neural network regression for continuous pain intensity estimation in video. *CoRR*, abs/1605.00894, 2016.
- [119] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. *CoRR*, abs/1711.06396, 2017.
- [120] M.Zeeshan Zia, Michael Stark, and Konrad Schindler. Are cars just 3d boxes? - jointly estimating the 3d shape of multiple objects. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.

- [121] Zeeshan Zia, Michael Stark, Bernt Schiele, and Konrad Schindler. Detailed 3d representations for object recognition and modeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 35(11):2608 – 2623, 2013.