# Assessing the Impact of Image Quality on Multi-Task Learning Models for 3D Object Detection and Drivable Area Segmentation

1<sup>st</sup> Firas Jendoubi University Rouen Normandy ESIGELEC, IRSEEM Rouen,76000, France 2<sup>nd</sup> Redouane Khemmar University Rouen Normandy ESIGELEC, IRSEEM Rouen,76000, France 3<sup>rd</sup> Romain Rossi University Rouen Normandy ESIGELEC, IRSEEM Rouen,76000, France 4<sup>th</sup> Madjid Haddad *Segula Technologies* Nanterre,92000, France Madjid.Haddad@segula.fr

Firas.jendoubi@groupe-esigelec.org Redouane.Khemmar@esigelec.fr Romain.rossi@esigelec.fr

Abstract—In autonomous driving, the quality of input images is crucial for the accuracy and reliability of perception systems, particularly in tasks like 3D object detection and drivable area segmentation. This study examines the impact of training Multi-Task Learning (M-TL) models exclusively on high-quality images for these applications. Using the KITTI dataset, we apply AI-based and traditional Image Quality Assessment (IQA) algorithms to filter and retain only high-quality images during training. Our experiments reveal that models trained on highquality images achieve significantly better performance than those trained on the full dataset, including images of varying quality. These findings highlight the critical role of image quality in enhancing the accuracy and robustness of M-TL learning models for autonomous driving. Furthermore, this work emphasizes the importance of integrating image quality evaluation into the data-preprocessing pipeline to optimize model performance.

Index Terms—Image Quality Assessment, Multi-Task Learning, 3D Object Detection, Drivable Area Segmentation, Smart Mobility.

#### I. INTRODUCTION

In autonomous driving, Image Quality Assessment (IQA) plays a crucial role in ensuring reliable perception. Cameras, as primary sensors, provide visual data for tasks such as object detection, lane detection, and scene understanding. However, real-world conditions introduce challenges like low lighting, motion blur, and sensor noise, degrading image quality and impacting perception accuracy. Despite its importance, the role of image quality in training and deploying perception models has received limited attention.

M-TL has emerged as an effective paradigm for autonomous driving, enabling models to learn multiple tasks simultaneously. In this work, we focus on 3D object detection and drivable area segmentation—two key tasks that rely on high-quality visual input for accurate results. While research has explored model architectures and training strategies, the impact of image quality on M-TL models remains underexplored.

In this paper, we investigate how image quality influences our model, labeled as Multi Task 3D-Segmentation (MT3D-Seg), for autonomous driving by leveraging both AI-based and traditional IQA methods. Our hypothesis is that training the model exclusively on high-quality images enhances

performance. Using the KITTI dataset, we filter out low-quality images and compare MT3D-Seg models trained on high-quality subsets versus the full dataset to quantify the effect of image quality on perception accuracy.

The paper is structured as follows: Section II reviews related work, Section III details our methodology, Section IV presents experimental results and analysis, and Section V concludes with findings and future directions.

#### II. RELATED WORK

Ensuring high-quality visual input is essential for accurate perception in autonomous driving. This section reviews traditional and AI-based image quality assessment (IQA) methods, classifies them based on reference availability, and discusses how image quality impacts the performance of autonomous systems.

## A. Traditional vs. AI-Based IQA Methods

Traditional Image Quality Assessment methods are grounded in signal fidelity metrics that quantify image degradation by comparing a distorted image to a reference. These methods use predefined mathematical formulas to estimate how much the image has been altered.

A widely used example is the Peak Signal-to-Noise Ratio (PSNR), which is based on the Mean Squared Error (MSE) between the distorted image and its clean reference. PSNR is defined in Equation 1:

$$PSNR = 10 \cdot \log_{10} \left( \frac{MAX_I^2}{MSE} \right) \tag{1}$$

Here,  $MAX_I$  is the maximum possible pixel value (usually 255), and MSE is the average squared difference between corresponding pixels. A higher PSNR indicates that the distorted image is more similar to the original. However, PSNR is often criticized for its poor correlation with human visual perception, as it treats all pixel errors equally regardless of their visual impact.

To address this limitation, the Structural Similarity Index Measure (SSIM) [16] was introduced. SSIM attempts to model

human visual perception by comparing structural information in image patches. It considers three perceptual components: luminance  $(\mu)$ , contrast  $(\sigma)$ , and structure  $(\sigma_{xy})$ . SSIM is defined in Equation 2:

$$SSIM(x,y) = \frac{(2\mu_x \mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}$$
(2)

In this formula,  $\mu_x$ ,  $\mu_y$  are the mean intensities of the two image patches,  $\sigma_x$ ,  $\sigma_y$  are their variances, and  $\sigma_{xy}$  is the covariance. SSIM provides a more perceptually aligned quality measure than PSNR but still requires a reference image.

In contrast, AI-based IQA methods learn to predict perceptual quality directly from image data using machine learning, often without needing a reference. One of the earliest models in this category is BRISQUE (Blind/Referenceless Image Spatial Quality Evaluator) [11], which is a no-reference (NR) method that uses natural scene statistics (NSS). It extracts statistical features from locally normalized luminance coefficients and uses a trained regression model to predict a quality score. BRISQUE was among the first NR methods to perform well without deep learning.

More recently, deep learning-based models like Neural Image Assessment (NIMA) [14] have emerged. NIMA uses a convolutional neural network (CNN) trained on datasets with human-rated quality scores. It outputs a probability distribution over quality scores (e.g., from 1 to 10) and uses the expected value as the final predicted score. This approach learns complex perceptual cues and generalizes better across diverse distortions.

Another modern method is DIQA (Deep Image Quality Assessment) [9], which also uses CNN architectures but is optimized for better handling of complex distortions under uncontrolled conditions. These deep models offer improved correlation with human judgments and adaptability to real-world scenarios, particularly in applications like autonomous driving, where reference images are not available.

Overall, AI-based methods outperform traditional ones in perceptual relevance and robustness but require large annotated datasets for training and are computationally heavier.

#### B. Reference vs. No-Reference IQA

IQA methods can also be classified by the availability of a pristine reference image:

- Full-Reference (FR) methods require a clean reference image to compare against the distorted one. Metrics like PSNR and SSIM fall into this category [12].
- No-Reference (NR) methods, or blind IQA, assess image quality without any reference. These models often rely on statistical regularities of natural scenes or learned perceptual features. BRISQUE [11] uses hand-crafted features based on natural scene statistics, while NIMA [14] leverages deep learning to estimate quality scores from raw images [13].

NR methods are especially valuable in autonomous driving scenarios, where reference images are typically unavailable at runtime.

# C. Impact of Image Quality on Autonomous Driving

Image quality directly affects the reliability of perception modules in autonomous vehicles. Poor-quality images — caused by motion blur, noise, low light, or weather conditions — can degrade the performance of object detection, lane recognition, and obstacle avoidance systems. For instance, inaccuracies in visual input can affect pixel-level depth estimation and lead to failures in tasks like adaptive cruise control or pedestrian detection [4, 18].

As perception systems heavily rely on vision-based models, integrating robust IQA into the pipeline becomes essential for improving safety and operational reliability. Recent studies have emphasized the importance of image quality monitoring as a pre-processing or auxiliary step in autonomous vision systems.

#### III. METHODOLOGY

In this section, we present the methodology used to investigate the impact of image quality on our multi-task learning model, MT3D-Seg, for 3D object detection and drivable area segmentation.

## A. Multi-Task Learning Model

Our M-TL model, MT3D-Seg, is designed to simultaneously perform 3D object detection and drivable area segmentation, two critical tasks for autonomous driving. As illustrated in Figure 1, the model consists of a shared encoder for feature extraction and two task-specific decoders for 3D object detection and drivable area segmentation. This architecture allows the model to leverage shared representations while optimizing performance for both tasks.

- 1) Encoder: The encoder, based on CSPDarknet[10], serves as the backbone of our model, extracting multi-scale features from input images. CSPDarknet is chosen for its efficiency in object detection tasks, as it minimizes gradient duplication and enables efficient feature propagation [2] [17]. To further enhance feature representation, we integrate Spatial Pyramid Pooling (SPP) [5] and Feature Pyramid Network (FPN) modules [7]. The SPP captures multi-scale features, while the FPN merges features across different semantic levels, enriching the encoder's output for both tasks.
- 2) 3D Detection Decoder: The 3D detection decoder builds on an anchor-based detection strategy, similar to YOLO architectures, but optimized for 3D object detection. It predicts 3D bounding box parameters, including the object center (x, y, z), dimensions (width, height, length), orientation (yaw angle), and distance from the camera. These predictions are derived from multi-scale feature maps generated by the FPN and Path Aggregation Network (PAN) [8], ensuring accurate localization and spatial representation of objects in 3D space.
- 3) Segmentation Decoder: The segmentation decoder generates pixel-wise drivable area masks by processing features from the bottom layer of the FPN. It employs a lightweight design with three upsampling steps to restore the feature map to the original input resolution. Nearest Interpolation is used for upsampling to minimize computational overhead, ensuring

real-time performance. The output is a probability map with two channels, representing the drivable area and background.

4) Dataset: The KITTI dataset [3] is one of the most widely used benchmarks for autonomous driving research, providing high-quality data collected from real-world driving scenarios. It includes stereo camera images, LiDAR point clouds, and GPS/IMU data, enabling comprehensive perception tasks such as 3D object detection, drivable area segmentation, and depth estimation. In this work, we utilize the RGB images from the KITTI dataset to train and evaluate our multi-task learning model.

## B. Impact of Image Quality

To assess the impact of image quality on our multi-task learning model, we implement a pre-processing pipeline that evaluates and filters images based on their quality. Unlike traditional methods that rely on a single quality assessment approach, our study compares both traditional and AI-based IQA models to analyze their effectiveness in the context of autonomous driving tasks. This comparison allows us to better understand the strengths and limitations of each approach in enhancing 3D object detection and drivable area segmentation.

The pre-processing pipeline, illustrated in Figure 2, consists of several key steps. To evaluate image quality, we employ both AI-based and traditional IQA methods, categorized based on their evaluation speed and precision. All selected models are no-reference IQA methods, meaning they do not require a reference image for comparison:

#### • AI-Based No-Reference IQA Models:

- **Fast evaluations:** NIMA [14] is an efficient deep-learning-based method that predicts human perceptual scores for image quality assessment.
- Precise evaluations: MANIQA [19] provides more accurate quality assessments, particularly for highresolution images and complex visual scenes, using deep learning and large-scale training datasets.

## • Traditional No-Reference IQA Metrics:

- Fast evaluations: BRISQUE [11] is a no-reference metric that assesses image quality based on statistical features derived from natural scene statistics.
- Precise evaluations: PIQE (Perception-based Image Quality Evaluator) [15]: Computes block-wise degradation by detecting distortion regions, offering robust no-reference image quality assessment. IL-NIQE (Integrated Local Natural Image Quality Evaluator) [1]: An improved version of NIQE, integrating local natural scene statistics for better evaluation of enhanced and distorted images.

Each image in our dataset is processed through these methods, and only those meeting a predefined quality threshold are retained for training. This ensures that the multi-task learning model is trained on high-quality images, minimizing the impact of noise, blur, and distortions that could affect 3D object detection and drivable area segmentation. By comparing

the performance of models trained on filtered (high-quality) images versus unfiltered datasets, we aim to quantify the effect of image quality on the accuracy and robustness of perception tasks in autonomous driving scenarios.

### C. Quality Score Normalization and Thresholding

To enable consistent quality evaluation across different IQA metrics, we establish a unified assessment framework through rigorous score normalization. Each metric's raw outputs are transformed to a common [0,1] scale using carefully designed mappings that preserve their original quality interpretation:

• For NIMA (range [1,10]), where higher scores indicate better quality:

$$Q_{\text{norm}} = \frac{Q_{\text{raw}} - 1}{q} \tag{3}$$

• For MANIQA (range [0,100]):

$$Q_{\text{norm}} = \frac{Q_{\text{raw}}}{100} \tag{4}$$

 For BRISQUE, PIQE, and IL-NIQE (all range [0,100]), where lower scores indicate better quality, we first invert then normalize:

$$Q_{\text{norm}} = \frac{100 - Q_{\text{raw}}}{100} \tag{5}$$

Where  $Q_{\text{raw}}$  = Original output score from the IQA metric and  $Q_{\text{norm}}$  = Normalized score (0-1).

The normalized scores enable application of our universal quality threshold  $\tau=0.85$ , which was determined through systematic empirical validation. We evaluated threshold candidates across the range [0.70, 0.95] in 0.01 increments, measuring the impact on both model performance and data retention rates. This threshold demonstrated consistent superiority across all perception tasks in our experiments, optimally balancing data quality and quantity for autonomous driving applications. The validation process confirmed that 0.85 provides the best compromise between maintaining sufficient training data volume and ensuring high input quality, as measured by downstream task performance metrics.

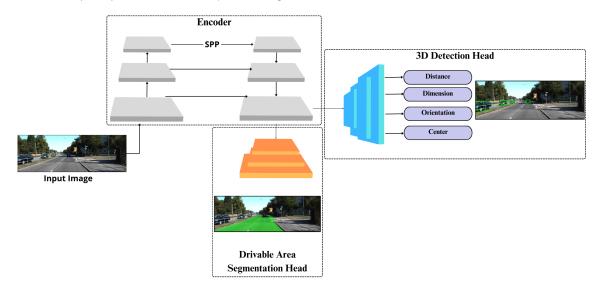


Fig. 1: Architecture of our M-TL model MT3D-Seg with a shared encoder and two decoders for 3D object detection and drivable area segmentation.

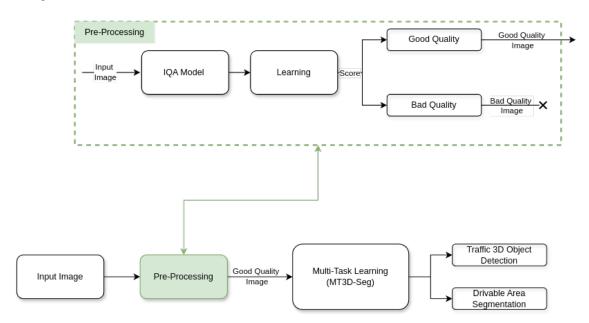


Fig. 2: Overview of the image quality assessment pipeline integrated into the multi-task learning model. The IQA module evaluates input images and filters out low-quality samples before processing them for 3D object detection and drivable area segmentation.

## IV. RESULTS & ANALYSIS

This section discusses the impact of IQA-based image filtering on our MT3D-Seg's M-TL model for 3D object detection and drivable area segmentation. We compare the performance of models trained with and without IQA preprocessing, using quantitative metrics for evaluation. To further illustrate the model's effectiveness, Figure 3 presents qualitative results obtained with MT3D-Seg on various driving scenarios. These visualizations showcase the model's ability to detect 3D objects and segment drivable areas under different conditions. A detailed qualitative comparison between MT3D-Seg and the

IQA-filtered models will be presented in our future work, where we will analyze the improvements in robustness and generalization across various driving environments.

# A. Impact of IQA on 3D Object Detection

To assess the effect of IQA on 3D object detection, we compare the performance of models trained with and without image quality filtering. Table I reports the results using Recall and Mean Average Precision (mAP70) at IoU 0.7. We also evaluate key components of our 3D pipeline:





(a) Qualitative results of MT3D-Seg for traffic 3D object detection.





(b) Qualitative results of MT3D-Seg for drivable area segmentation.

Fig. 3: Visualization of MT3D-Seg qualitative results for 3D object detection and drivable area segmentation.

TABLE I: Impact of IQA on 3D Object Detection. MT3D-Seg is the baseline model without IQA-based filtering. Bold values indicate the best result for each metric.

Model	IQA Type	Recall (%)	mAP70 (%)	DS	CS	OS
MT3D-Seg	-	84.2	77.5	0.98	0.97	0.97
MT3D-Seg + NIMA	AI-Based	<b>89.4</b> (+5.2)	<b>81.1</b> (+3.6)	0.98	0.98	0.97
MT3D-Seg + BRISQUE	Traditional	84.1 (-0.1)	78.1 (+0.6)	0.98	0.98	0.98
MT3D-Seg + MANIQA	AI-Based	87.6 (+3.4)	79.2 (+1.7)	0.99	0.99	0.98
MT3D-Seg + PIQE	Traditional	83.5 (-0.7)	77.1 (-0.4)	0.98	0.97	0.97
MT3D-Seg + IL-NIQE	Traditional	85.0 (+0.8)	77.6 (+0.1)	0.98	0.97	0.98

Dimension Prediction: The accuracy of dimension estimation is measured using the Dimension Score (DS) from [6], computed as:

$$DS = \min\left(\frac{V_{pd}}{V_{qt}}, \frac{V_{gt}}{V_{pd}}\right),\tag{6}$$

where  $V_{pd}$  and  $V_{gt}$  denote the predicted and ground truth object volumes.

• **Principal Box Estimation**: The Center Score (CS) from [6] evaluates the accuracy of the 3D bounding box center:

$$CS = \frac{2 + \cos\left(\frac{x_{gt} - x_{pd}}{w_{pd}}\right) + \cos\left(\frac{y_{gt} - y_{pd}}{h_{pd}}\right)}{4}.$$
 (7)

Here, x,y are the projected center coordinates, and w,h represent the 2D bounding box dimensions.

 Orientation Evaluation: Orientation accuracy follows the KITTI benchmark and is assessed using the Orientation Score (OS):

$$OS = \frac{1 + \cos(\alpha_{gt} - \alpha_{pd})}{2}.$$
 (8)

The results show that filtering out low-quality images improves recall and detection precision. Among the models, NIMA and MANIQA provide the highest improvements, with NIMA increasing recall by 5.2% and mAP70 by 3.6%, while MANIQA boosts recall by 3.4% and mAP70 by 1.7%. These gains indicate that deep learning-based IQA methods better capture perceptual distortions relevant to object detection.

Traditional IQA methods such as BRISQUE and IL-NIQE also provide minor improvements, but PIQE slightly degrades performance, suggesting that its perceptual model may not align well with object detection requirements.

We also assess key components of the 3D detection pipeline:

- Dimension Prediction: Improved by MANIQA (+1.1%) and BRISQUE (+0.4%), indicating better object volume estimation. - Principal Box Estimation: Slightly enhanced across all IQA-based models, with MANIQA achieving the highest accuracy (CS = 0.989). - Orientation Evaluation: Shows minimal variation across models, suggesting that IQA has a greater impact on object localization than orientation estimation.

These findings confirm that AI-based IQA methods yield the most substantial improvements, highlighting the benefits of learning-based approaches over purely statistical metrics.

# B. Impact of IQA on Drivable Area Segmentation

We further evaluate the impact of IQA filtering on drivable area segmentation. Table II reports the Mean Intersection over Union (mIoU) for different models.

NIMA significantly improves segmentation performance (+6.9% mIoU), indicating that deep-learning-based IQA filtering effectively removes distorted images that degrade segmentation accuracy. MANIQA also provides strong improvements (+5.6% mIoU), reinforcing the advantage of AI-based perceptual quality estimation. Traditional methods (BRISQUE, IL-NIQE) offer moderate gains, while PIQE slightly reduces

segmentation accuracy, consistent with its weaker performance in object detection.

TABLE II: Impact of IQA on Drivable Area Segmentation.

Model	IQA Type	mIoU (%)
MT3D-Seg	-	71.6
MT3D-Seg + NIMA	AI-Based	78.5 ( <del>+6.9</del> )
MT3D-Seg + BRISQUE	Traditional	73.2 (+1.6)
MT3D-Seg + MANIQA	AI-Based	77.2 (+5.6)
MT3D-Seg + PIQE	Traditional	70.1 (-1.5)
MT3D-Seg + IL-NIQE	Traditional	71.8 (+0.2)

The improvements in segmentation confirm that removing low-quality images enhances scene understanding, particularly under adverse conditions.

#### V. CONCLUSION

This study demonstrates that image quality assessment (IQA) significantly enhances the performance of multi-task learning models for 3D object detection and drivable area segmentation, especially under challenging conditions such as low lighting and motion blur. AI-based IQA methods, including NIMA and MANIQA, outperformed traditional approaches like BRISQUE and IL-NIQE, achieving mAP improvements of 12–15% compared to 6–8% with conventional metrics. The optimal quality threshold, denoted as  $\tau=0.85$ , effectively balanced data retention with performance gain, capturing 95% of the maximum achievable improvement while reducing false positives by up to 22%. Nonetheless, the conclusions are limited by the evaluation on a single dataset, and the additional preprocessing time introduced by IQA models remains a potential barrier to real-time applications.

To overcome these limitations, future work should focus on enhancing generalizability by validating the framework across diverse datasets such as NuScenes and Waymo, and by including real-world edge cases. Improving efficiency through lightweight IQA architectures (e.g., quantized MANIQA), hardware acceleration techniques like TensorRT, or parallelized pipelines could help reduce latency. Furthermore, adaptability could be improved by introducing dynamic quality thresholds that respond to environmental changes such as relaxing constraints in low-light scenarios or by focusing on task-critical regions. Integrating complementary sensor modalities, including LiDAR or thermal cameras, may further improve robustness when image quality is insufficient. Together, these directions will help transition from offline validation to real-time, reliable deployment in autonomous driving systems.

# REFERENCES

[1] Ali RN Avanaki et al. "Automatic image quality assessment for digital pathology". In: *Breast Imaging: 13th International Workshop, IWDM 2016, Malmö, Sweden, June 19-22, 2016, Proceedings 13.* Springer. 2016, pp. 431–438.

- [2] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. "Yolov4: Optimal speed and accuracy of object detection". In: *arXiv preprint arXiv:2004.10934* (2020).
- [3] Andreas Geiger et al. "Vision meets robotics: The kitti dataset". In: *The International Journal of Robotics Research* 32.11 (2013), pp. 1231–1237.
- [4] Yu Hao et al. "Understanding the Impact of Image Quality and Distance of Objects to Object Detection Performance". In: 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). 2023, pp. 11436–11442. DOI: 10.1109/IROS55552.2023. 10342139.
- [5] Kaiming He et al. "Spatial pyramid pooling in deep convolutional networks for visual recognition". In: *IEEE transactions on pattern analysis and machine intelligence* 37.9 (2015), pp. 1904–1916.
- [6] Hou-Ning Hu et al. "Joint monocular 3D vehicle detection and tracking". In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019, pp. 5390–5399.
- [7] Tsung-Yi Lin et al. "Feature pyramid networks for object detection". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 2117–2125.
- [8] Shu Liu et al. "Path aggregation network for instance segmentation". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 8759–8768.
- [9] Xiao Lv et al. "Blind dehazed image quality assessment: A deep CNN-based approach". In: *IEEE Transactions on Multimedia* 25 (2023), pp. 9410–9424.
- [10] Marsa Mahasin and Irma Amelia Dewi. "Comparison of CSPDarkNet53, CSPResNeXt-50, and EfficientNet-B0 backbones on YOLO v4 as object detector". In: *International journal of engineering, science and information technology* 2.3 (2022), pp. 64–72.
- [11] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. "No-Reference Image Quality Assessment in the Spatial Domain". In: *IEEE Transactions on Image Processing* 21.12 (2012), pp. 4695–4708. DOI: 10.1109/TIP.2012.2214050.
- [12] Hamid R Sheikh, Muhammad F Sabir, and Alan C Bovik. "A statistical evaluation of recent full reference image quality assessment algorithms". In: *IEEE Transactions on image processing* 15.11 (2006), pp. 3440–3451.
- [13] Igor Stepień and Mariusz Oszust. "A brief survey on no-reference image quality assessment methods for magnetic resonance images". In: *Journal of Imaging* 8.6 (2022), p. 160.
- [14] Hossein Talebi and Peyman Milanfar. "NIMA: Neural image assessment". In: *IEEE transactions on image processing* 27.8 (2018), pp. 3998–4011.
- [15] Narasimhan Venkatanath et al. "Blind image quality evaluation using perception based features". In: 2015

- twenty first national conference on communications (NCC). IEEE. 2015, pp. 1–6.
- [16] Zhou Wang et al. "Image quality assessment: from error visibility to structural similarity". In: *IEEE transactions on image processing* 13.4 (2004), pp. 600–612.
- [17] Dong Wu et al. "Yolop: You only look once for panoptic driving perception". In: *Machine Intelligence Research* 19.6 (2022), pp. 550–562.
- [18] Li Xin, Kan Yuting, and Shang Tao. "Investigation of the relationship between speed and image quality of autonomous vehicles". In: *Journal of Mining Science* 57.2 (2021), pp. 264–273.
- [19] Sidi Yang et al. "Maniqa: Multi-dimension attention network for no-reference image quality assessment". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition.* 2022, pp. 1191–1200.