Active Data Acquisition in Autonomous Driving Simulation

Jianyu Lai, Zexuan Jia, Boao Li

Abstract—Autonomous driving algorithms rely heavily on learning-based models, which require large datasets for training. However, there is often a large amount of redundant information in these datasets, while collecting and processing these datasets can be time-consuming and expensive. To address this issue, this paper proposes the concept of an active data-collecting strategy. For high-quality data, increasing the collection density can improve the overall quality of the dataset, ultimately achieving similar or even better results than the original dataset with lower labeling costs and smaller dataset sizes. In this paper, we design experiments to verify the quality of the collected dataset and to demonstrate this strategy can significantly reduce labeling costs and dataset size while improving the overall quality of the dataset, leading to better performance of autonomous driving systems. The source code implementing the proposed approach is publicly available on the page 1.

Index Terms—Autonomous driving, Simulation, Data Acquisition, Active Strategy

I. INTRODUCTION

Currently, most of the algorithms related to autonomous driving use learning-based models, which typically require a large amount of data as support. Therefore, the size of autonomous driving datasets is often extremely large (TB level). Collecting these data and using them to train autonomous driving algorithms incurs high costs in terms of both collection time and training time. However, in practice, there is often a large amount of redundant information in the largescale datasets used for autonomous driving. This redundant information can only have a small positive impact on the training of algorithms, but the cost of collecting and using this redundant information for training is very high. Therefore, our goal is to develop a new data collection tool that can reduce the redundant information in the collected dataset and improve the quality of the dataset, thus increasing the training efficiency of the algorithm and reducing the time and cost of data collection and training by providing the data to the algorithm before training.

In order to improve the quality of the dataset, the precision of data points and the quantity of data are considered as two metrics. It is vital to find a balance between them: for high-quality data, increasing the collection density can improve the overall quality of the dataset, ultimately achieving similar or even better results than the original dataset with lower labelling costs and smaller dataset size. In this project, all experiments were completed and tested in a simulated

Jianyu Lai, Zexuan Jia, Boao Li are with the Department of Computer Science and Technology, Southern University of Science and Technology, Shenzhen, 518055, China

scene. First, it is cost-effective. Conducting tests in real-world scenarios can be time-consuming and expensive [1], [2], [?]. Using simulation scenes can significantly reduce the cost of testing. By creating virtual scenarios, researchers can quickly test various driving scenarios and abnormal situations, which can improve testing efficiency. Second, it is convenient to do data recording and analysis. In simulation scenes, various data such as sensor data and driving trajectories can be easily recorded and analyzed. These data can be used to improve the algorithms and performance of automatic-driving systems.

Vehicle detection is an important research topic in the field of autonomous driving, and one of the hotspots in the field of computer vision [3]. Because of the strong dependency of the object detection algorithm, the quality of a dataset must be evaluated by a specific algorithm, thus one kind of algorithm must me selected as the criterion for the evaluation of the dataset [4], [5]. For autonomous vehicles, the comprehension of the surroundings is of primary importance, and more data and images lead to better performance. However, considering the project scale and the time usage, it is necessary to choose a proper algorithm that is not only time-saved but also can be trained well with images at mid-scale. Considering the difficulty and the engineering requirement, the detection algorithm YOLO v5 is chosen. This series of algorithms will be introduced in the following section.

II. RELATED WORK

This section only briefly discusses our literal review of several popular object detection algorithms. Since object detection requires object localization and object type classification [6], the learned features can be categorized into one-stage and two-stage object detection algorithms [7].

One-stage object detection algorithms generate object location and object classification results directly in one stage, so they do not require a region proposal process, which is usually simpler and faster than two-stage detection algorithms [8]. There are two kinds of one-stage detection algorithms, single-shot multi-box detector (SSD) [9] and you only look once (YOLO) [10]. The SSD algorithm uses a VGG16 network as the backbone, and SSD-512 as the input version. While the Yolo detection algorithm is a real-time object detection algorithm, and the Yolo-v3 is the most popular version because of the model's compact size and ease of application [11], which uses the DarkNet-53 as the backbone.

The two-stage object detection algorithms need to conduct a region proposal process first and then classify the object in the proposed region [4], [5]. The basic algorithm is the R-CNN algorithm. However, the two-stage construction means

¹https://github.com/Th1nkMore/carla_dataset_tools

more time will be used in the training processing, so the R-CNN algorithm is not a light-weighted algorithm so two advanced algorithms, mask R-CNN and faster R-CNN will be introduced. The Fast R-CNN uses DCNN(deep CNN) as the backbone, and it employs a region proposal algorithm that reduces the number of object region proposals and improves the region proposal quality to improve the learning speed [12]. The Mask R-CNN enhances the overall detection accuracy and small target detection, which the pool layers affect a lot [13].

Measure	Faster R-CNN	YOLOv3
TP (True Positives)	578	751
FP (False Positives)	2	2
FN (False Negatives)	150	7
Precision (TPR)	99.66%	99.73%
Sensitivity (recall)	79.40%	99.07%
F1 Score	88.38%	99.94%
Quality	79.17%	98.81%
Processing Time (Av. in ms)	1.39s	0.057ms

TABLE I: Faster R-CNN and YOLO v3

From Table I [14], the Yolo-v3 performs better than the fastest two-stage detection algorithm faster R-CNN in all attributes, which reduces the FN a lot and increase the quality and processing time, which means the FPS(frame per second) is much higher and the mAP is more accurate. And From Table II, all the Yolo series algorithms and faster R-CNN are compared. We can find that the mAP of Yolo-v5 and Yolo-v3 is similar, but Yolo-v5 has about 6 times FPS of Yolo-v3 and 2 times FPS as Yolo-v4, which shows that it is very efficient. Because in our project, the scale of the dataset and the training time both count, so after consideration, the light-weighted and accurate algorithm Yolo-v5 is chosen.

Model	Size (pixels)	Test dataset	mAP (0.5)	FPS	GPU
Fast R-CNN	600×1000	VOC2007	70	0.5	Titan X
Faster R-CNN	600×1000	VOC2007	76.4	5	Titan X
YOLOv1	446×448	VOC2007	63.4	45	Titan X
SSD	512×512	VOC2007	76.8	19	Titan X
YOLOv2	544×544	MS COCO	44.0	40	Titan X
YOLOv3	608×608	MS COCO	57.9	20	Titan X
YOLOv4	608×608	MS COCO	65.7	62	Tesla V100
YOLOv5s	640×640	MS COCO	55.4	113	Tesla V100

TABLE II: The specific comparison among YOLO series algorithms.

For the simulation environment, the CARLA [15] simulator is chosen. CARLA is an open-source simulator for autonomous driving research. CARLA has been developed from the ground up to support the development, training, and validation of autonomous urban driving systems. In addition to open-source code and protocols, CARLA provides open digital assets (urban layouts, buildings, vehicles) that were created for this purpose and can be used freely. The simulation platform supports flexible specifications of sensor suites and environmental conditions [16], [17], [18]. CARLA is generally used to study the performance of three approaches to autonomous driving: a classic modular pipeline, an endto-end model trained via imitation learning, and an end-toend model trained via reinforcement learning. The approaches are evaluated in controlled scenarios of increasing difficulty, and their performance is examined via metrics provided by CARLA, illustrating the platform's utility for autonomous

driving research. In addition, the motion flow method can be easily used in the data collected in the CARLA simulation environment, which can assist in the analysis of the trajectory of the cars. In CARLA, the target ODD can be specified more easily than in real environments. In this paper, all the experiments are conducted in and tested in a CARLA simulation environment.

III. METHODOLOGY

A. Indicator algorithms

In object detection tasks, the quality and quantity of training data are critical factors that affect the performance of the algorithm. For autonomous driving cars, understanding the surrounding environment is the most important thing [19], the YOLO algorithm which is chosen, the more instance is provided, the better the performance is. While increasing the number of instances can lead to better performance, it is not always feasible to collect or annotate a large amount of data. In this work, we propose a novel approach that utilizes the information provided by YOLO to generate synthetic data that can be used to improve the performance of the algorithm. By analyzing the confidence scores and bounding boxes output by YOLO, we can identify the most informative instances and use them to generate high-quality synthetic data. In this section, we present our methodology for data generation and evaluation, and demonstrate the effectiveness of our approach on a benchmark dataset.

B. Data Complexity Design

What is more, to test the algorithm efficiently, the data collection method and strategies are critical. For a specific algorithm, the diversity and the quality will affect the training result. For diversity, we use the UQI (Universal Image Quality Index) [20] as one of the metrics. When the similarity of two consecutive images is higher than one threshold, the image will be removed. This method can grant the diversity of the training dataset and prove that no similar images will be used so that the training time will be reduced as well. The formula is as follows:

$$Q = \frac{4\sigma_{xy}\bar{x}\bar{y}}{(\sigma_x^2 + \sigma_y^2)[(\bar{x}^2 + \bar{y}^2)]}$$

When it comes to quality, there are several factors, but the main two factors that we consider are the number of objects in each image and the occlusion degree of each object in one image. For the first one, when one image has more target objects, the YOLO algorithm can be trained better with the same data scale, because more features can be shown and learned by the algorithm. Then, for the occlusion degree, because the semantic segmentation image classifies all the same kind of objects in the same color, so when two or more objects are near each other, they might be labelled as one big car. So reducing or recognizing this kind of image are very critical.

C. Radar Algorithms

Radar cross-section (RCS) measurements [21] require a well-calibrated radar system. Often, the calibration process is performed by evaluating the received echo signal from a target with a well-known RCS. In a ground-based setup, radar measurements may be significantly affected by the environment in general and by the multipath in particular. Also, Driver assistance and safety systems are in high demand by customers and legislators. Those systems are working satisfactorily on highways, but for urban scenarios, a more precise position, orientation and dimension estimation is essential for time-crucial systems. It is assumed that the target vehicle does not drift, and hence the orientation can be used to determine the direction of movement [22].

D. Experiment Methods

For the experiment and testing aspect, the main input of this project is the semantic segmentation of images captured through cameras in the Carla simulation environment, along with instance segmentation to assist in judgment and detection. Semantic segmentation can automatically classify different types of objects in the scene, facilitating annotation tools to annotate them. Instance segmentation is easy to observe and understand, making it easy to check the results. In the data of the Carla simulation tool, all cars and traffic lights (objects to be detected) have their own ground truth, i.e., their actual positions in the image. The method used in this article is to first use existing labelling tools to annotate the semantic segmentation images, and then use the annotated data to train the YOLO algorithm. Finally, give the trained YOLO algorithm some raw segmentation data, let it label the detected objects, and then check the results of properties such as IoU [23] and mAP50.

Regarding image quality, there are several influencing factors, but the two main factors are the number of objects in each image and the degree of occlusion of each object in each image. For the first one, when an image has more target objects, using the same data size can better train the YOLO algorithm because the algorithm can display and learn more features. As for occlusion, since the semantic segmentation image classifies all objects of the same type as the same colour when two or more objects are close to each other, they may be labeled as larger objects. Reducing or identifying such images is critical.

IV. EXPERIMENTAL SETUP

In order to evaluate the effectiveness of the active data collection strategy proposed in this work, which may bring a higher unit data value, two experiments were designed and conducted. The first experiment aimed to compare the performance of a model trained on data collected using the active strategy with the performance of a model trained on data collected using a passive strategy while controlling for other variables. The second experiment aimed to compare the performance of a model trained on a fixed amount of data collected using the active strategy with the performance of a

model trained on the same amount of data collected using a passive strategy, while also controlling for other variables.

Here is a brief view of the dataset configuration, which also shows the experiment design:

- D1: Collect all data passively
- D1-S: Active collect for the same time as D1
- D2:Active collect until the same size as D1
- V: Dataset used for validate and test

Dataset name	Total Frames	Map
D1	900	Town02
D1-S	579	Town02
D2	900	Town02
V	375	Town03

TABLE III: Dataset Configuration

A. Equivalent Time-samples Experiment

This experiment was conducted by collecting two datasets, denoted as dataset **D1** and dataset **D1-S** respectively. Dataset **D1-S** was collected using the active data collection strategy proposed in this work, while dataset **D1** collects all data. The datasets were collected at the same time, with the same map and vehicle number, in order to make the experiment comparative.

The performance of a model trained on dataset **D1** was then compared with the performance of a model trained on dataset **D1-S**. Specifically, the metrics of the models trained on the two datasets were compared, and the time required to train the models was also recorded. We expect that the metrics of the model trained on dataset **D1-S** will be similar to the metrics of the model trained on dataset **D1**, but the time required to train the model using dataset **D1-S** will be significantly less than the time required to train the model using dataset **D1**.

B. Equivalent Size-samples Experiment

This experiment was conducted by collecting two datasets, denoted as dataset **D1** and dataset **D2** respectively. Data set **D2** uses the active strategy and data set **D1** not uses the strategy respectively, and changes map and vehicle number to make the experiment comparative, then observes the time and effect of the training model. The datasets were collected at the same size, with the same map and vehicle number, in order to make the experiment comparative.

The performance of a model trained on dataset **D1** was then compared with the performance of a model trained on dataset **D2**. Specifically, the metrics of the models trained on the two datasets were compared, and the time required to train the models was also recorded. We expect that the metrics of the model trained on dataset **D2** will be significantly better than the metrics of the model trained on dataset **D1**, and that the time required to train the model using dataset **D1** will be similar to the time required to train the model using dataset **D2**.

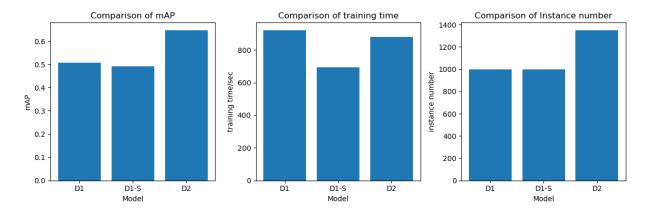


Fig. 1: Comparison of mAP, training time and instance number

Dataset name	Vehicle	Training time	mAP	mAP50 95
D1	996	14min40s	0.508	0.391
D1-S	996	11min32s	0.492	0.365
D2	1835	16min07s	0.647	0.525

TABLE IV: Raw Result

V. RESULTS

From the equivalent Time-samples Experiment, we observed that by employing an active acquisition strategy, the number of collected images decreased while the number of instances remained unchanged. The effectiveness of the approach remained similar. However, the training time decreased from 14 minutes and 40 seconds to 11 minutes and 32 seconds (about **21.36%** reduction).

From the equivalent Size-samples Experiment, under the same amount of data, the model trained by the dataset obtained using an active acquisition strategy demonstrated a significant increase in the number of instances collected. The training time remained comparable, but there was a noticeable improvement in performance (approximately **27.36**% increment).

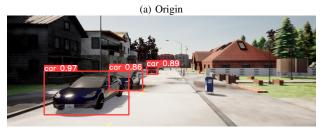
When dealing with images that contain multiple instances, the performance of a model trained using the active collecting strategy is generally better. This is primarily because the model has learned from a dataset that includes a diverse range of images with multiple instances within a single image.

Through this exposure, the model becomes more adept at detecting, localizing, and understanding multiple instances within a complex image. It learns to recognize patterns, relationships, and spatial configurations between different objects, enabling it to accurately identify and differentiate between the individual instances present.

This enhanced understanding of multiple instances leads to improved performance when the model is applied to real-world scenarios or tasks involving images with multiple objects. The model's ability to generalize from the diverse training examples it encountered during the active collecting process helps it better handle instances that may vary in appearance, pose, scale, or occlusion.

Furthermore, the active collecting strategy allows the model to focus on areas that are challenging or underrepresented







(c) Passive collecting

Fig. 2: Comparison of models. Note that the confidences are different.

in the dataset. This targeted selection helps address potential biases or limitations present in the initial training dataset, ensuring that the model receives sufficient exposure to instances with different characteristics, thereby reducing the risk of overfitting.

In summary, training a model with an active collecting strategy proves advantageous when dealing with images containing multiple instances. The model's exposure to diverse examples and its improved ability to understand and differentiate between objects within a single image contribute to its superior performance in tasks related to multi-instance image analysis.

In summary, the experimental results are generally consistent with our expectations. The outcomes align with the an-

ticipated trends and patterns we hypothesized. These findings support our initial assumptions and provide further evidence for the effectiveness and viability of the proposed methods. The observed outcomes validate the hypotheses formulated and encourage future investigations. Additional studies could explore various factors and variables to enhance the robustness and generalizability of the findings.

VI. CONCLUSIONS

In this paper, we proposed an active data acquisition strategy for improving the performance of the algorithm of autonomous driving, which is YOLO v5. We evaluated the effectiveness of this strategy through two experiments, which demonstrated that the data collected using this strategy has a higher unit data value compared to the data collected using a passive strategy. Specifically, we found that the models trained on data collected using the active strategy achieved higher accuracy and required less time to train compared to the models trained on data collected using a passive strategy. These results suggest that the active data acquisition strategy proposed in this work can be used to improve the performance of autonomous driving algorithms in various aspects. By collecting informative data points and reducing the data acquisition time, the proposed strategy can help to reduce the cost and time required to train autonomous driving algorithm models, while maintaining or improving the accuracy of the algorithms. In summary, the experiments have demonstrated that the hypothesis proposed in this paper is correct. For the vehicle detection problem, improving the Ground Truth generation tool will enable the proposed method to be used to verify and evaluate various active parameters and corresponding model performance, thereby validating the quality of the dataset. Active autonomous driving data collection, from the perspective of the dataset, may become another path for further improving autonomous driving algorithms. In the future, we will extend this work with other AI technologies, such as knowledge graph [24], [25] and human-centered AI [26] to reduce the reality gap [27], [28].

ACKNOWLEDGEMENTS

REFERENCES

- G. Lan, J. M. Tomczak, D. M. Roijers, and A. Eiben, "Time efficiency in optimization with a bayesian-evolutionary algorithm," *Swarm and Evolutionary Computation*, vol. 69, p. 100970, 2022.
- [2] G. Lan, M. De Carlo, F. van Diggelen, J. M. Tomczak, D. M. Roijers, and A. E. Eiben, "Learning directed locomotion in modular robots with evolvable morphologies," *Applied Soft Computing*, vol. 111, p. 107688, 2021.
- [3] X. Dong, S. Yan, and C. Duan, "A lightweight vehicles detection network model based on yolov5," in *Engineering Applications of Artificial Intelligence*, vol. 113, 2022, p. 104914.
- [4] G. Lan, L. De Vries, and S. Wang, "Evolving efficient deep neural networks for real-time object recognition," in 2019 IEEE Symposium Series on Computational Intelligence (SSCI). IEEE, 2019, pp. 2571– 2578.
- [5] G. Lan, J. Benito-Picazo, D. M. Roijers, E. Domínguez, and A. Eiben, "Real-time robot vision on low-performance computing hardware," in 2018 15th international conference on control, automation, robotics and vision (ICARCV). IEEE, 2018, pp. 1959–1965.
- [6] Z. Gao and G. Lan, "A neat-based multiclass classification method with class binarization," in *Proceedings of the genetic and evolutionary* computation conference companion, 2021, pp. 277–278.

- [7] Y. Xiao, Z. Tian, J. Yu, Y. Zhang, S. Liu, S. Du, and X. Lan, "A review of object detection based on deep learning," *Multimedia Tools and Applications*, vol. 79, no. 33, pp. 23729–23791, 2020.
- [8] G. Lan, Z. Gao, L. Tong, and T. Liu, "Class binarization to neuroevolution for multiclass classification," *Neural Computing and Applications*, vol. 34, no. 22, pp. 19845–19862, 2022.
- [9] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [10] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE* conference on computer vision and pattern recognition, 2016, pp. 779– 788.
- [11] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," CoRR, vol. abs/1804.02767, 2018.
- [12] R. Girshick, "Fast r-cnn," in Proceedings of the IEEE international conference on computer vision, 2015, pp. 1440–1448.
- [13] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in Proceedings of the IEEE international conference on computer vision, 2017, pp. 2961–2969.
- [14] B. Benjdira, T. Khursheed, A. Koubaa, A. Ammar, and K. Ouni, "Car detection using unmanned aerial vehicles: Comparison between faster r-cnn and yolov3," in 2019 1st International Conference on Unmanned Vehicle Systems-Oman (UVS). IEEE, 2019, pp. 1–6.
- [15] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proceedings of the 1st Annual Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, S. Levine, V. Vanhoucke, and K. Goldberg, Eds., vol. 78. PMLR, 13–15 Nov 2017, pp. 1–16.
- [16] G. Lan, Z. Luo, and Q. Hao, "Development of a virtual reality teleconference system using distributed depth sensors," in 2016 2nd IEEE International Conference on Computer and Communications (ICCC). IEEE, 2016, pp. 975–978.
- [17] G. Lan, J. Sun, C. Li, Z. Ou, Z. Luo, J. Liang, and Q. Hao, "Development of uav based virtual reality systems," in 2016 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI). IEEE, 2016, pp. 481–486.
- [18] T. Xiang, F. Jiang, G. Lan, J. Sun, G. Liu, Q. Hao, and C. Wang, "Uav based target tracking and recognition," in 2016 IEEE international conference on multisensor fusion and integration for intelligent systems (MFI). IEEE, 2016, pp. 400–405.
- [19] H. Xu, G. Lan, S. Wu, and Q. Hao, "Online intelligent calibration of cameras and lidars for autonomous driving systems," in 2019 IEEE Intelligent Transportation Systems Conference (ITSC). IEEE, 2019, pp. 3913–3920.
- [20] Z. Wang and A. Bovik, "A universal image quality index," *IEEE Signal Processing Letters*, vol. 9, no. 3, pp. 81–84, 2002.
- [21] G. Galati, G. Pavan, and C. Wasserzier, "Environmental effects on ground-based radar measurements," in 2019 IEEE 5th International Workshop on Metrology for AeroSpace (MetroAeroSpace), 2019, pp. 349–354.
- [22] F. Roos, D. Kellner, J. Klappstein, J. Dickmann, K. Dietmayer, K. D. Muller-Glaser, and C. Waldschmidt, "Estimation of the orientation of vehicles in high-resolution radar images," in 2015 IEEE MTT-S International Conference on Microwaves for Intelligent Mobility (ICMIM), 2015, pp. 1–4.
- [23] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. Huang, "UnitBox," in Proceedings of the 24th ACM international conference on Multimedia. ACM, oct 2016.
- [24] T. Liu, G. Lan, K. A. Feenstra, Z. Huang, and J. Heringa, "Towards a knowledge graph for pre-/probiotics and microbiota-gut-brain axis diseases," *Scientific Reports*, vol. 12, no. 1, p. 18977, 2022.
- [25] G. Lan, T. Liu, X. Wang, X. Pan, and Z. Huang, "A semantic web technology index," *Scientific reports*, vol. 12, no. 1, p. 3672, 2022.
- [26] G. Lan, Y. Wu, F. Hu, and Q. Hao, "Vision-based human pose estimation via deep learning: A survey," *IEEE Transactions on Human-Machine* Systems, 2022.
- [27] G. Lan, J. Chen, and A. Eiben, "Evolutionary predator-prey robot systems: From simulation to real world," in *Proceedings of the genetic* and evolutionary computation conference companion, 2019, pp. 123– 124.
- [28] G. Lan, J. Chen, and A. E. Eiben, "Simulated and real-world evolution of predator robots," in 2019 IEEE Symposium Series on Computational Intelligence (SSCI). IEEE, 2019, pp. 1974–1981.