

Article

Not peer-reviewed version

Benchmarking Perception to Streaming Inputs in Vision-Centric Autonomous Driving

[Tianshi Jin](#) , [Weiping Ding](#) ^{*} , Mingliang Yang , [Honglin Zhu](#) , Peisong Dai

Posted Date: 21 November 2023

doi: 10.20944/preprints202311.1184.v1

Keywords: vision-centric perception benchmark; online assessment; streaming inputs; two-dimensional entropy



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

Benchmarking Perception to Streaming Inputs in Vision-Centric Autonomous Driving

Tianshi Jin, Weiping Ding *, Mingliang Yang, Honglin Zhu and Peisong Dai

¹ School of Mechanical Engineering, Southwest Jiaotong University, Chengdu 610031, China; ts.jin@my.swjtu.edu.cn (T.J.);

* Correspondence: dwp@swjtu.edu.cn

Abstract: In recent years, vision-centric perception has played a crucial role in autonomous driving tasks, encompassing functions such as 3D detection, map construction, and motion forecasting. However, the deployment of vision-centric approaches in practical scenarios is hindered by substantial latency, often deviating significantly from the outcomes achieved through offline training. This disparity arises from the fact that conventional benchmarks for autonomous driving perception predominantly conduct offline evaluations, thereby largely overlooking the latency concerns prevalent in real-world deployment. While a few benchmarks have been proposed to address this limitation by introducing effective evaluation methods for online perception, they do not adequately consider the intricacies introduced by the complexity of input information streams. To address this gap, we propose the Autonomous-driving Streaming I/O (ASIO) benchmark, aiming to assess the streaming inputs characteristics and online performance of vision-centric perception in autonomous driving. To facilitate this evaluation across diverse streaming inputs, we initially establish a dataset based on the CARLA Leaderboard. In alignment with real-world deployment considerations, we further develop evaluation metrics based on information complexity specifically tailored for streaming inputs and streaming performance. Experimental results indicate significant variations in model performance and ranking under different major camera deployments, underscoring the necessity of thoroughly accounting for the influences of model latency and streaming inputs characteristics during real-world deployment. To enhance streaming performance consistently across distinct streaming inputs features, we introduce a backbone switcher based on the identified streaming inputs characteristics. Experimental validation demonstrates its efficacy in perpetually improving streaming performance across varying streaming inputs features.

Keywords: vision-centric perception benchmark; online assessment; streaming inputs; two-dimensional entropy

1. Introduction

Vision-centric perception has attracted considerable attention in the field of autonomous driving in recent years. It is intuitive that vision plays the most dominant role in human driving. In principle, vision-centric perception can obtain the richest semantic information, which is essential for decision-making in autonomous driving, compared to LiDAR-based and millimeter-wave radar-based perception. Moreover, we found a large body of previous research on vision-based perception for various autonomous driving tasks in past years, such as 3D detection [1–10] in driving scenes, map construction [11–15], motion prediction [16,17], and even end-to-end autonomous driving [18–20].

Despite the remarkable achievements in vision-centric perception research, many methods suffer from high latency when deployed in real-world settings, which hinders their online performance. For example, in 3D detection, a fundamental task for autonomous driving, camera-based 3D detectors usually have much longer (See Table 1) inference time than LiDAR-based counterparts [21–23] (on NVIDIA RTX4090). Therefore, it is essential to have evaluation metrics that balance accuracy and latency. However, most of the existing benchmarks [24–34] focus on evaluating

offline performance only (e.g., Average Precision (AP), Intersection over Union (IoU), etc.). Although some studies have adopted the streaming Perception paradigm to measure accuracy-latency trade-offs, and Wang et al. [35] proposed an online evaluation protocol that can assess the online performance of different perception methods under various hardware conditions, they still lack prior evaluation of the streaming input. This means that for the online performance evaluation of the vision-centric perception, they are still missing the initial impact of the streaming inputs.

This paper introduces the Autonomous-driving Streaming I/O (ASIO) benchmark to address the problems mentioned above. The benchmark quantifies the effects of different input sources on streaming perception performance and fills the gap in existing research. Current research only focuses on evaluating online performance without considering the influence of various streaming inputs on the perception system (e.g., different resolutions, field of view angles, etc.). Unlike mainstream datasets and benchmarks, our benchmark is based on the CARLA Leaderboard [36] simulator, ensuring consistency in the environment and targets encountered during testing. We collect real-time streaming inputs data and manually annotate the targets to measure the online performance of 3D detection. Practical deployment is also investigated, specifically the problem of ASIO under different inputs. We design evaluation metrics for perception streaming inputs based on the fractional dimensional entropy computation method of time series to assess streaming perception with different models. Figure 1 illustrates the significant impact of perception inputs' variation on the streaming performance of different methods. Our approach provides a more precise characterization of the effect of perception input on the deployment of real-world autonomous driving tasks compared to classical offline benchmarks. The main contributions of this paper are summarized as follows:

- (1) We present the ASIO benchmark, for quantitatively evaluating the characteristics of camera-based streaming input and the streaming perception performance, which opens up possibilities for vision-centric perception design and performance prediction for autonomous driving.
- (2) A scenario and dataset for evaluating different streaming inputs are built based on the CARLA Leaderboard, which enables camera-based 3D detection streaming evaluation.
- (3) For the implicit characteristics in streaming inputs, the computation of fractional order entropy values in one and two dimensions is proposed to construct quantitative metrics, where we investigate the streaming performance of seven modern camera-based 3D detectors under various streaming inputs.

The remainder of this paper is organized as follows: in "Section 2," we analyze current offline and online evaluation methods for autonomous driving perception and identify research gaps within them. Additionally, we delve into methods for characterizing information complexity. Subsequently, in "Section 3," we introduce the dataset established for assessing streaming inputs, along with the establishment of metrics for online evaluation. Meanwhile, we propose an improved method for 3D detector enhancement based on the aforementioned approaches. In "Section 4," we conduct a streaming performance evaluation of seven typical 3D detectors across various cameras, extracting features of streaming perception and validating the a priori nature of our metrics. This section identifies factors that should be considered in the practical deployment of perception systems. "Section 5" concludes the entire work, highlighting existing issues and providing suggestions for future improvements.

Table 1. Comparison between autonomous-driving perception datasets.

Dataset	Stream.	Modality	Task	Model speed
KITTI [26]	×	L&C	Multi-task	-
Argoverse [32]	×	L&C	Multi-task	-
nuScenes [24]	×	L&C	Multi-task	-
Waymo [31]	×	L&C	Multi-task	-
CARLA Leaderboard	×	L&C	Multi-task	-
Argoverse-HD [37]	√	C	2D det.	~40FPS
nuScenes-H [35]	√	C	3D det.	~7FPS

Waymo	√	L	3D det.	~25FPS
CARLA Leaderboard	√	C	2D&3D det.	~30FPS@2D det. ~5FPS@3D det.

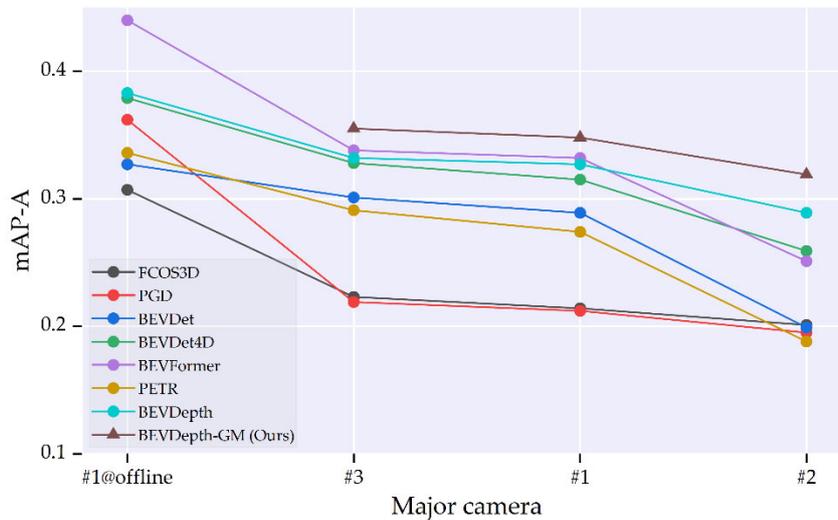


Figure 1. Comparison of streaming performances on our benchmark, where the model rank changes on streaming inputs' variation. And the detector BEVDepth-GM (built on BEVDepth [4]) equipped with our switcher achieved better streaming performance on different major cameras.

2. Related Work

2.1. Autonomous-Driving Benchmark

Thanks to the various open-source benchmarks, the past decade has witnessed significant progress in autonomous driving perception. These benchmarks have shifted from 2D detection [25,28,29,33,34] tasks to 3D detection [24,26,27,30–32] tasks, which are tailored for autonomous driving scene understanding. Additionally, the data acquisition has also progressed from single RGB images to multi-modal data. Even popular datasets with 3D annotations employ surround-view images, greatly facilitating the development of vision-centric perception. However, these benchmarks primarily focus on assessing the offline performance of perception algorithms, overlooking the practical issues of perception system deployment

2.2. Streaming Perception

The deployment of perception in autonomous driving faces the challenge of balancing accuracy and latency. In order to improve perception performance, previous works have explored the concept of streaming perception, which utilizes temporal information. For example, Li et al. [37] introduced a benchmark for image detection algorithms and proposed a method based on Kalman filtering [38] and reinforcement learning [39] to mitigate latency. Han et al. [40] developed an efficient streaming detector for LiDAR-based 3D detection tasks, accurately predicting future frames. Wang et al. [35] presented a benchmark for various perception models under different computational constraints, with a focus on the 3D detection task. These works establish evaluation paradigms for camera-based 3D detection, as well as LiDAR-based 3D detection, highlighting the trade-off between accuracy and latency in real-world deployment. However, it is also worth noting the significant impact of streaming input on overall performance [35,37]. Therefore, there is a need for a methodology that incorporates input sources into the streaming perception evaluation for autonomous driving.

2.3. Nonlinear Time Series Complexity

Autonomous driving perception systems are similar to and based on real-world e.g., ecological, meteorological, geological, etc. systems generated by natural or physical mechanisms, and are complex systems that are difficult to explain their modes of operation deterministically or by constructing analytical models [41,42]. When dealing with streaming input, it is important to quantify its complexity. The concept of Shannon's information entropy has been continuously promoted and extended, leading to the proposal of various discrete forms of entropy metrics such as Rényi entropy [43], Tsallis entropy [44], approximate entropy [45], sample entropy [46], and permutation entropy [47]. These metrics have become the main tools for measuring the complexity of a system. Ubriaco [48] introduced a new entropy measure known as fractional entropy. This measure promotes the integer order Shannon entropy in fractional dimensions, providing the possibility of entropy metrics in the application of systems with long-range correlation. The fractional entropy not only has a high sensitivity to the dynamic changes of the signal features but also reveals more details and information about the system. As a result, it demonstrates better utility in practice [49]. To explore the implicit information, we incorporated this analysis into the ASIO benchmark. We considered one- and two-dimensional aspects in the processing of the streaming input, aiming to reveal the underlying information.

3. Methods

This section begins with an introduction to the concept of ASIO. Then, we provide a test scenario and corresponding dataset to evaluate the holistic streaming perception. Finally, we present evaluation metrics to measure the streaming perception performance across various input conditions.

3.1. Autonomous-Driving Streaming I/O

The evaluation of the ASIO benchmark has two aspects: evaluating the information complexity of the streaming inputs and the streaming perception performance online.

Obviously, autonomous driving perception is a multiple-input, multiple-output system at each level, and usually, the dimensions of the inputs are larger than the dimensions of the outputs, and each level can be represented by the model in Figure 2. The input vector X consists of components X_1, X_2, \dots, X_m , and the output vector Y consists of components Y_1, Y_2, \dots, Y_l , so that the requirement of the subsystem is to maximize the information transmitted to the output about the inputs of the system, according to the principle of maximum mutual information. Based on the information-theoretic model described above, it is necessary to evaluate the complexity of the streaming inputs X . Specially, given streaming inputs $\{X_m\}_{m=1}^T$, where is the image inputs at timestamp t_m and T is the total number of input timestamps. The perception algorithms are acquired to make an online response to the input instance, and the entire online predictions are $\{\hat{Y}_l\}_{l=1}^U$, where \hat{Y} is the prediction at the timestamp t_l , and U represents the total number of predictions. Notably, the prediction timestamps are not synchronized with the input timestamps, and the model inference speed is typically slower than the input frame rate (i.e., $U < T$). To evaluate the predictions at the input timestamp t_m , the ground truth Y_m is desired to match with the most recent prediction, yielding the pair $(Y_m, \hat{Y}_{\theta(m)})$, where $\theta(m) = \arg \max_l t_l < t_m$. Based on the matching strategy, the ASIO benchmark evaluates the complexity of the streaming inputs:

$$\mathcal{H}_{\text{ASIO}} = \mathcal{H}(X_m), \quad (1)$$

and the online performance at every input timestamp:

$$\mathcal{O}_{\text{ASIO}} = \mathcal{O}\left(\{(Y_m, \hat{Y}_{\theta(m)})\}_{m=1}^T\right), \quad (2)$$

where $\mathcal{H}(\cdot)$ and $\mathcal{O}(\cdot)$ are the evaluation metrics, which will be elaborated in subsequent sections. Notably, ASIO instantiates the streaming paradigm on camera-based 3D detection, and the key insights also generalize to other vision-centric perceptions in autonomous driving.

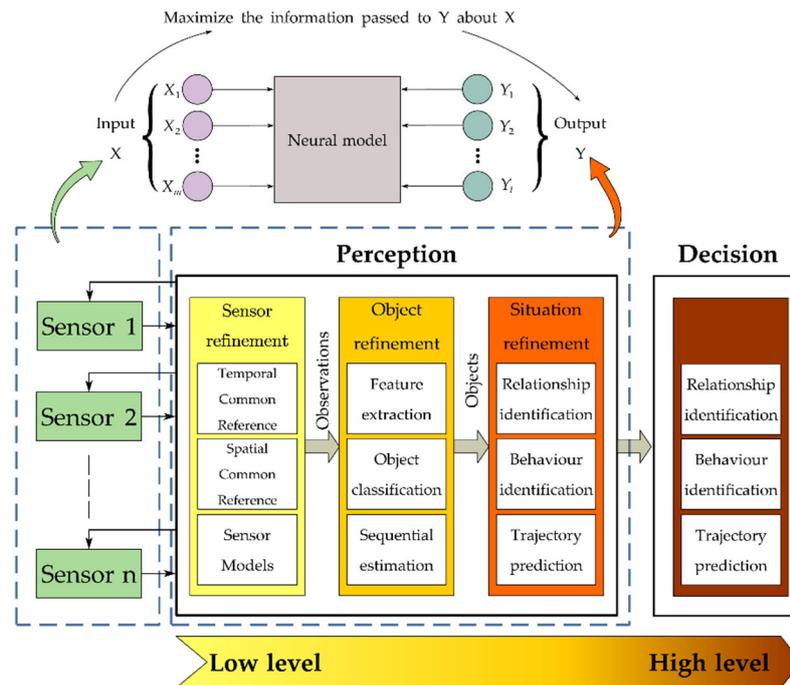


Figure 2. The perception system process encompasses various levels of data, information processing, and modeling in an overall scheme. We view the sensor's inputs and perceived performance in the model as the I/O of streaming.

3.2. Scenarios and Dataset

To validate all the proposed methods, we established a standardized visual perception benchmark using Carla. The evaluation scenario map Town 12 from the CARLA Leaderboard 2 was chosen for this purpose (Figure 3). This city consists of various contrasting regions, including urban, residential, and rural areas, with a surrounding highway system and a ring road. The architectural styles in the city resemble those commonly seen in medium to large cities worldwide. In our study, we focused primarily on urban scenes and selected a fixed number of road participants, such as vehicles of different types and colors. The proportions of these vehicles were determined based on a relevant survey. Additionally, the traffic participants and the ratios of cars to trucks and vehicles to pedestrians were designed accordingly. The scenario's roads encompassed various types of lanes, road markings, traffic lights, signage, and even different weather conditions (Table 2). In this scenario, the vehicle will travel along a fixed route at a fixed speed for a total distance of 5km. It will be equipped with a visual perception sensor that needs evaluation. The setup traffic participants, which are equipped with their respective motion states, will be traversed by the vehicle during the entire journey.

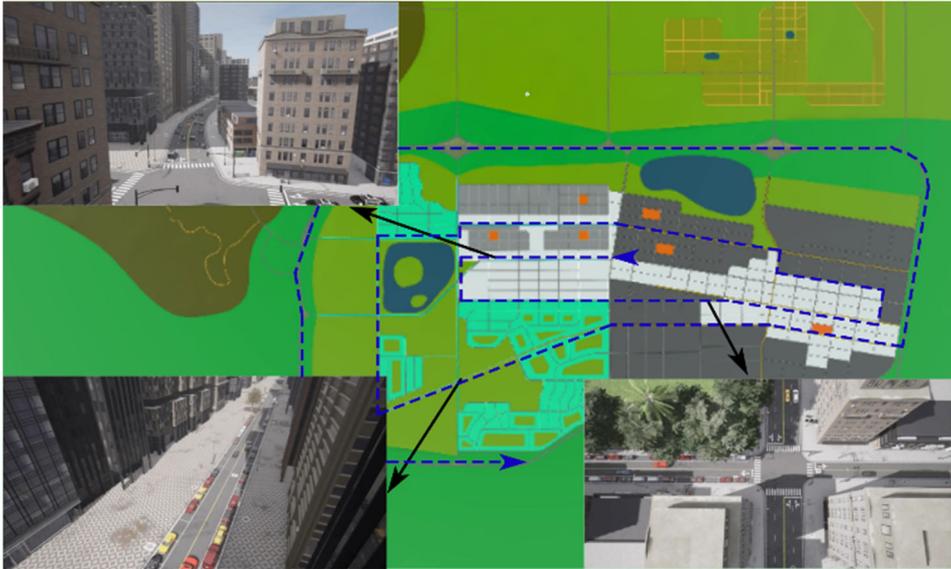


Figure 3. The simulation scenarios.

Table 2. Traffic participants, weather.

Types	Quantities/Discriptions	
Vehicles	200	
Pedestrians	30	
Traffic Lights	Default in map	
Sidewalks	Default in map	
Weather	Noon	Clear
		Mid Rainy
	Night	Clear
		Mid Rainy

Meanwhile, we set up an RGB baseline camera with a sampling frequency of 10Hz and a resolution of 1600×900 , and a 32-beam spinning LiDAR with a sampling frequency of 10Hz, 360° horizontal FOV, and -30° to 10° vertical FOV for annotation purposes. For each weather condition, we annotated the data at 10 Hz for a cycle of 1000s, resulting in 10000 images and point cloud frames to be annotated. We conducted tests under four weather conditions while keeping the route and the targets constant, so the annotation work did not need to be repeated. For 3D detection, we adopted the SUSTechPoints annotation tool. For each streaming input to be tested, we combined the annotated point cloud information and used the CAT [53] model to complete the 3D box annotation and motion state assignment on the test images. In contrast to the existing streaming perception datasets (see Table 3), we built a simulator-based dataset by manually collecting and annotating data, enabling the evaluation of streaming inputs. We tested some perception models on our constructed dataset and compared them with the mainstream dataset nuScenes (see Table 4). In light of constraints pertaining to the dataset capacity and the limited variety and quantity of target objects we incorporated, the outcomes reveal an overestimation of scores within our dataset. Nevertheless, the overall performance exhibits trends analogous to those observed in nuScenes, thereby providing partial validation of the validity of our dataset.

Table 3. Streaming perception dataset comparison.

Dataset	Construction Methods	Task	Evaluation Containing		
			Inputs	Computation	Online
Argoverse-HD2	1. Based on Argoverse	2D Detection	×	×	√
	2. Extended with manually added annotations				
nuScenes -H	1. Based on nuScenes	3D Detection	×	√	√
	2. Expanding the annotations therein from 3 Hz to 12 Hz.				
Ours	1. Based on CARLA	3D Detection	√	×	√
	2. Manually annotation of baseline				
	3. Automatic annotation of test object				

Table 4. Popular algorithm validation on our dataset.

Methods	nuScenes		Ours	
	mAP	NDS	mAP	NDS
FCOS3D [50]	0.358	0.428	0.459	0.538
DETR3D [51]	0.412	0.479	0.520	0.580
BEVFormer [5]	0.481	0.569	0.581	0.688
BEVDepth	0.520	0.609	0.629	0.719
SOLOFusion [10]	0.540	0.619	0.647	0.721

3.3. Evaluation Metrics

We developed evaluation metrics for streaming inputs and streaming performance, aiming to comprehensively examine the holistic streaming perception of various 3D detectors under different inputs. This section first introduces the streaming inputs metrics and then explains the streaming performance metrics.

3.3.1. Streaming Inputs Metrics.

As shown in the information-theoretic model in Figure 2, we need to reveal the information complexity of the streaming inputs sequence. One common way to evaluate the pixel inputs is to calculate their information entropy. To describe the local structural features of the streaming inputs, we introduce two-dimensional entropy which reveals the combined features of pixel grey scale information and grey scale distribution in the vicinity of the pixel. The feature pair (x_1, x_2) is formed by the gray level of the current pixel and the mean value of its neighborhood. Here, x_1 represents the gray level of the pixel and x_2 represents the mean value of the neighbors. The combined probability density distribution function of x_1 and x_2 is then given by the following equation:

$$p(x_1, x_2) = \frac{f(x_1, x_2)}{P \times Q}, \quad (3)$$

where $f(x_1, x_2)$ is the frequency at which the feature pair (x_1, x_2) appears, and the size of X is $P \times Q$. In our implementation, x_1 is derived from the eight adjacent neighbors of the center pixel, as depicted in Figure 4. The discrete fractional two-dimensional entropy is defined as follows:

$$H = - \sum_{x_1=0}^{255} \sum_{x_2=0}^{255} p(x_1, x_2) \log_2 p(x_1, x_2). \quad (4)$$

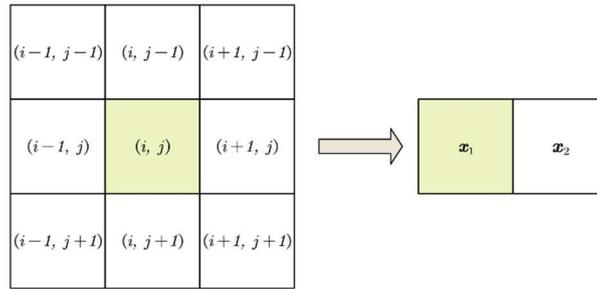


Figure 4. The pair (x_1, x_2) - a pixel and its eight neighborhoods

To achieve a long-range correlation of information metrics in perception systems, we use a fractional-dimensional extension of information entropy and generalize it. First, based on the theory of fractional calculus, the integer-order Shannon entropy is extended to fractional dimension, leading to the proposition of a new entropy measure called fractional entropy, as:

$$S_q(P) = \sum_i^n p_i (-\log_2 p_i)^q, \quad (5)$$

Proposed as an alternative measure to Shannon's information entropy, Cumulative Residual Entropy (CRE) [54] was introduced due to the difficulty in estimating the differential entropy of a continuous random variable through empirical distribution in practice. For a nonnegative discrete random variable X , its CRE is defined as:

$$\mathcal{E}(X) = -\sum \bar{F}(x) \log_2 \bar{F}(x), \quad (6)$$

where $F(x)$ is the cumulative distribution function of X and $\bar{F}(x) = 1 - F(x)$, it can be estimated through the empirical entropy value of the sample.

The single frames H computed above are represented as fractional dimensional CREs to obtain our proposed metrics for evaluating the information complexity of streaming inputs:

$$\mathcal{E}^q(X) = \sum \bar{F}(x) [-\log_2 \bar{F}(x)]^q, q \in [0,1]. \quad (7)$$

$$\mathcal{H} = \mathcal{E}^q(H). \quad (8)$$

The computed \mathcal{H} exhibit notable numerical disparities when dealing with sequential inputs containing pixels of varying field of view (FOV). In order to more accurately assess streaming inputs, we also incorporate the density of \mathcal{H} in FOV space, which we denote as \mathcal{D} :

$$\mathcal{D} = \frac{\mathcal{H}}{\log_2(H \times V)}, \quad (9)$$

where $H \times V$ represents the horizontal and vertical field of view angles.

Figure 5 shows the information complexity calculation of the streaming inputs from the baseline camera during a complete run of our benchmark. The run used a sliding window with overlap, with the window length set to 400 and the sliding distance set to 200, and calculated its \mathcal{H} values in fractional dimension order $q \in [0,1]$ of 0.02 steps. When key targets or scenarios occur, the values of \mathcal{H} at $q < 0.5$ show dramatic fluctuations with time scales. When q is close to 1, the values are smaller and vary slightly with time. Therefore, compared to \mathcal{H} , the classical approach (when $q = 1$) cannot effectively reveal the fluctuation of information from the streaming inputs about key targets and scenes. From this perspective, \mathcal{H} is superior to classical information entropy methods. by introducing fractional order parameters q , \mathcal{H} is able to capture the detailed variations of the system information, thus more accurately detecting the changes in the system's dynamic features and providing effective cues for evaluating the performance of visual perception. Numerous studies have shown that the introduction of the fractional dimension makes the entropy metric more applicable to the study of time series compared to Shannon entropy, which is also clearly demonstrated by the above results.

To obtain the q -value that is suitable for our benchmark, we expect the test \mathcal{H} to show minimal variation under different weather conditions, and to exhibit a larger gradient when critical scenarios occur. Therefore, the following two operators are proposed:

$$\sigma = \sqrt{\frac{\sum_{w1-w4} (\mathcal{H} - \bar{\mathcal{H}})^2}{n}}, \quad (10)$$

$$\nabla h = \sum_n \frac{h(s-d) - h(s+d)}{2d}, \quad (11)$$

where σ represents the standard deviation of \mathcal{H} in w1-w4 under the sliding total number n , \mathcal{H} is denoted by $h(s)$, and the value of d is the sliding window step size. By constructing the following optimization equation:

$$\begin{cases} \min f(\alpha) = [\sigma, -\nabla h] \\ \text{s.t. } q \in [0,1] \end{cases} \quad (12)$$

we obtain a $q = 0.36$ that is applicable to our benchmark under the current conditions, and this result is also used by default in the experimental section later on.

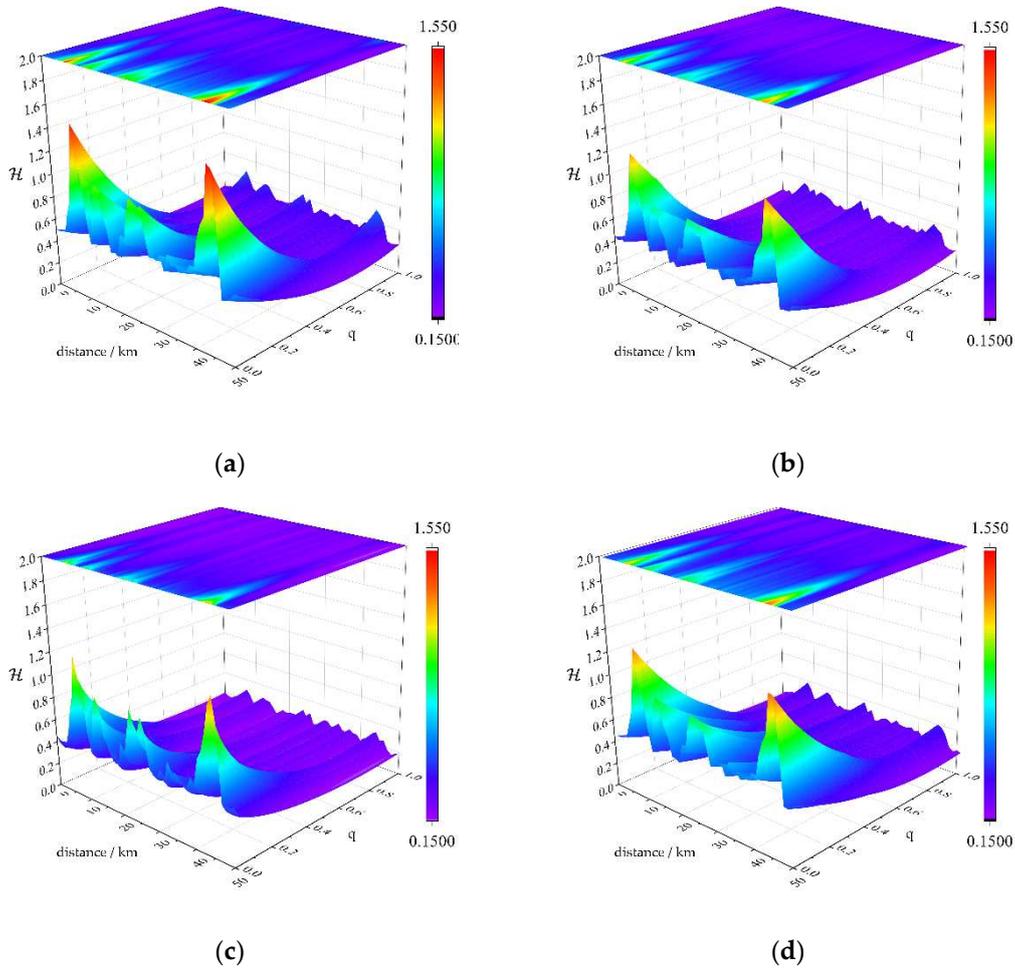


Figure 5. (a), (b), (c), and (d) depict the contour plots of the \mathcal{H} distributions in w1, w2, w3, w4, where w1 corresponds to Clear Noon, w2 to Rainy Noon, w3 to Clear Night, and w4 to Rainy Night. It is evident that, relative to the case where q is equal to 1 (Shannon's entropy description), the fractional dimensional order with q taking values within the range $(0, 1)$ provides a significantly more sensitive elucidation of the critical nodes in the benchmark scenario.

3.3.2. Streaming Performance Metrics.

We adopted the official nuScenes evaluation metrics, comprising Average Translation Error (ATE), Average Scale Error (ASE), Average Orientation Error (AOE), Average Velocity Error (AVE), Average Attribute Error (AAE), nuScenes Detection Score (NDS) and mean Average Precision (mAP). Due to the latency induced by the inference time in the streaming performance evaluation, there is a discrepancy between the predicted bounding boxes and the ground-truth locations, which is more pronounced for high-speed targets. Nevertheless, the AVE metric only quantifies the velocity error of true positive objects, which tend to be low-speed or static objects. Hence, we preserved the offline property of AVE, while applying streaming evaluation to other metrics, denoted as mAP-A, ATE-A, ASE-A, AOE-A, and AAE-A. Inspired by existing methods [35,37], we introduce ASIO Detection Score (ADS) as follows:

$$\text{ADS} = \frac{1}{10} [5\text{mAP-A} + \sum_{\text{mTP} \in \text{TP}} (1 - \min(1, \text{mTP}))], \quad (13)$$

where $\text{TP} = \{\text{AVE}, \text{ATE-A}, \text{ASE-A}, \text{AOE-A}, \text{AAE-A}\}$ is the set of true positive metrics.

3.3.3. ASIO Switcher

The ASIO benchmark evaluates the current detector inference at the instant time, regardless of its completion status. This way, the online inference latency will significantly affect the streaming perception performance. According to the conclusion of existing methods [35,37], the online performance is largely influenced by the streaming inputs and the backbone used by the detector, and it is very different from the offline evaluation result. It can be inferred that choosing the corresponding backbone for the 3D detector according to the complexity of different inputs can improve the online performance of the detector. Therefore, in this section, we design a switcher to different backbones for the detector by evaluating the information complexity of the current and previous streaming inputs.

For real-time sequential inputs, we select the information complexity measure $\mathcal{H}_i - \mathcal{H}_j$ within a time window, and observe that when the \mathcal{H} of streaming inputs fluctuates significantly, its local distribution approximates an exponential distribution. Therefore, we choose the grey model GM (1, 1) [55] to predict the \mathcal{H} values of the next k steps, and use the predicted values to construct a selector that decides whether to switch different backbone 3D detectors. The schematic diagram of this process is shown in Figure 6.

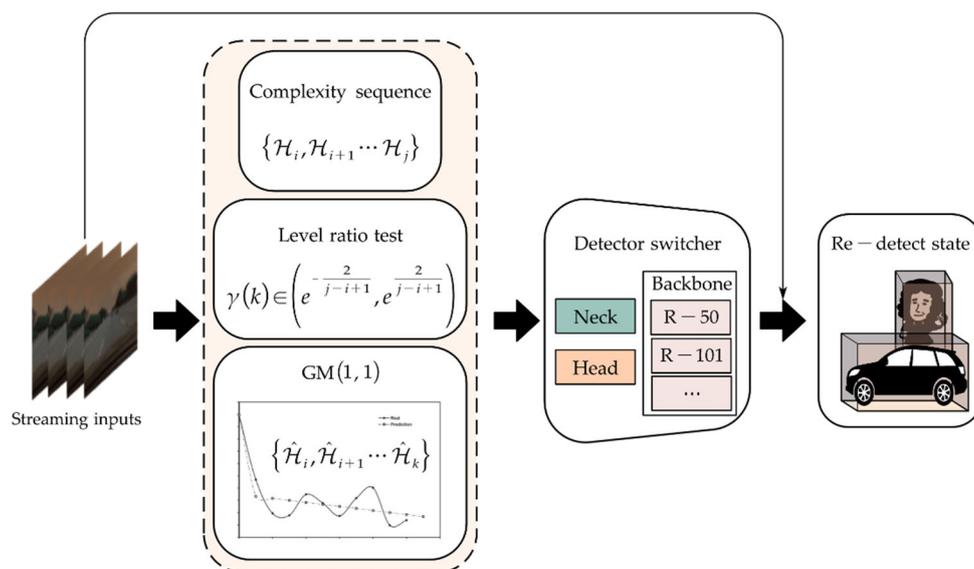


Figure 6. The streaming inputs updating switcher, where the GM (1, 1) is utilized to predict the future complexity of streaming inputs and is used as a basis for switching the backbone structure in the detector switcher and updating the re-detect state.

Assuming that we have obtained a sequence $\mathbb{H}^{(0)} = (\mathcal{H}^{(0)}(1), \mathcal{H}^{(0)}(2), \dots, \mathcal{H}^{(0)}(n))$, we test whether its ratio

$$\delta(k) = \frac{\mathcal{H}^{(0)}(k-1)}{\mathcal{H}^{(0)}(k)}, \quad k = 2, 3, \dots, n,$$

satisfies condition

$$\delta(k) \in \left(e^{-\frac{2}{n+1}}, e^{\frac{2}{n+1}} \right).$$

If this condition is satisfied, the grey model GM (1, 1) can be introduced at the present time. Let $\mathbb{H}^{(1)}$ be the 1-AGO sequence of $\mathbb{H}^{(0)}$,

$$\mathbb{H}^{(1)} = (\mathcal{H}^{(1)}(1), \mathcal{H}^{(1)}(2), \dots, \mathcal{H}^{(1)}(n)), \quad (14)$$

$$\mathcal{H}^{(1)}(k) = \sum_{i=1}^k \mathcal{H}^{(0)}(i), \quad k = 2, 3, \dots, n. \quad (15)$$

Additionally, obtain the sequence $\mathbb{Z}^{(1)}$ as the mean generating sequence of $\mathbb{H}^{(1)}$ within its immediate neighborhood:

$$\mathbb{Z}^{(1)} = (\mathcal{Z}^{(1)}(2), \mathcal{Z}^{(1)}(3), \dots, \mathcal{Z}^{(1)}(n)), \quad (16)$$

$$\mathcal{Z}^{(1)}(k) = 0.5\mathcal{H}^{(1)}(k) + 0.5\mathcal{H}^{(1)}(k-1). \quad (17)$$

We employ the elementary GM (1, 1) differential equation paradigm

$$\mathcal{H}^{(0)}(k) + a\mathcal{Z}^{(1)}(k) = b, \quad (18)$$

where a is called the development coefficient; parameters a and b are identified. Assume that:

$$\theta = \begin{bmatrix} a \\ b \end{bmatrix}^T, \quad Y_n = \begin{bmatrix} \mathcal{H}^{(0)}(2) \\ \mathcal{H}^{(0)}(3) \\ \vdots \\ \mathcal{H}^{(0)}(n) \end{bmatrix}, \quad B_n = \begin{bmatrix} -\mathcal{Z}^{(1)}(2) & 1 \\ -\mathcal{Z}^{(1)}(3) & 1 \\ \vdots & \vdots \\ -\mathcal{Z}^{(1)}(n) & 1 \end{bmatrix}$$

then matrix form can be written as:

$$Y_n = B_n \theta. \quad (19)$$

The objective function of least square method is established as:

$$J(\hat{\theta}_n) = (Y - B_n \hat{\theta}_n)^T (Y - B_n \hat{\theta}_n). \quad (20)$$

The least square method is used to get the estimated value $\hat{\theta}_n$ of $\theta = [a, b]^T$:

$$\hat{\theta}_n = (B_n^T B_n)^{-1} B_n^T Y_n. \quad (21)$$

We can get the whitening function of GM (1, 1) model:

$$\frac{d\mathcal{H}^{(1)}}{dt} + a\mathcal{H}^{(1)} = b. \quad (22)$$

The time response function of the available Equation (20) under initial conditions is:

$$\mathcal{H}^{(1)}(k) = \mathcal{H}^{(1)}(1) - bae^{-ak} + ba. \quad (23)$$

1-IAGO operation is applied to Equation (7), and the simulated sequence is obtained:

$$\mathcal{H}^{(0)}(k) = \mathcal{H}^{(1)}(k) - \mathcal{H}^{(1)}(k-1), \quad k = 2, 3, \dots, n. \quad (24)$$

By using the above prediction methods, we can obtain a series of predicted values when \mathcal{H} shows large fluctuations. These predicted values predict the information complexity of the streaming inputs, which we use as a basis for switching the detection algorithms for different backbones (V2-99, R50 and R101 were selected for the next experiments). Notably, this switcher can be applied to

any modern 3D detector (e.g., in the experiment, BEVDepth-GM is built upon BEVDepth). Moreover, the switching pipeline is lightweight (<10ms on CPU) and has negligible effect on streaming latency.

4. Experiments

In this section, the experiment setup is first given. Afterward, we explored evaluating different levels of input complexity, such as streaming inputs and online performance evaluation. Finally, we conducted real-vehicle experiments to confirm the accuracy of our benchmark.

4.1. Experiment Setup

In our benchmark, on the one hand inputs are subjected to complexity analysis and on the other hand vision-centric perception is instantiated as camera-based 3D detection. The dataset in section 3.2 is used to evaluate the inputs as well as the 3D detection. We followed the methods [35,37] and used a hardware simulator for the streaming evaluation. For the camera-based open-source detector to be tested, we measured its open-source code on a fixed GPU (NVIDIA RTX4090) with a batch size of 1.

4.2. Benchmarking Results

Using the ASIO benchmark, we analyzed seven modern autonomous driving 3D detectors (FCOS3D [50], PGD [54], BEVDet [1], BEVDet4D [2], BEVFormer [5], PETR [6], BEVDepth [4]) and our proposed BEVDepth-GM with three kinds of inputs from major sensors (see Table 5). As shown in Table 6, we can observe that:

Table 5. Major sensors for the three input types used for the experiments.

Camera	FOV/° H×V	Resolution	Frame rate
#1	90×65	1280×720	25
#2	123×116	848×800	25
#3	70×55	1024×768	25

Table 6. Comparison of the results of different autonomous driving 3D detectors on our dataset validation set, where BEVDepth-GM is based on our switcher built on BEVDepth. For different major cameras, we use the metrics in Section 3.3.1. For Streaming=×, we use the offline metrics. For Streaming=√, we use the online metrics in Section 3.3.2.

Methods	Major Camera	FPS	GLOPs	Streaming	\mathcal{H}	$\mathcal{D}(10x)$	ADS(NDS) ↑	mAP(-A) ↑
FCOS3D	#1	-	2008.2	×	-	-	0.376	0.307
	#1	1.5	2008.2	√	0.221	0.177	0.329	0.214
	#2	3.2	2008.2	√	0.210	0.152	0.318	0.201
	#3	1.7	2008.2	√	0.218	0.183	0.333	0.223
PGD	#1	-	2223.0	×	-	-	0.408	0.362
	#1	1.5	2223.0	√	0.221	0.177	0.337	0.212
	#2	3.0	2223.0	√	0.210	0.152	0.315	0.195
	#3	1.9	2223.0	√	0.218	0.183	0.347	0.219
BEVDet	#1	-	215.3	×	-	-	0.415	0.327
	#1	17.9	215.3	√	0.221	0.177	0.403	0.289
	#2	24.9	215.3	√	0.210	0.152	0.346	0.199
	#3	21.2	215.3	√	0.218	0.183	0.407	0.301
BEVDet4D	#1	-	222.0	×	-	-	0.480	0.379
	#1	14.5	222.0	√	0.221	0.177	0.445	0.315
	#2	23.3	222.0	√	0.210	0.152	0.383	0.259
	#3	16.3	222.0	√	0.218	0.183	0.455	0.328

BEVFormer	#1	-	1322.2	×	-	-	0.517	0.44
	#1	2.1	1322.2	√	0.221	0.177	0.461	0.332
	#2	5.0	1322.2	√	0.210	0.152	0.392	0.251
	#3	3.6	1322.2	√	0.218	0.183	0.475	0.338
PETR	#1	-	297.2	×	-	-	0.371	0.336
	#1	5.9	297.2	√	0.221	0.177	0.351	0.274
	#2	13.7	297.2	√	0.210	0.152	0.283	0.188
	#3	7.4	297.2	√	0.218	0.183	0.355	0.291
BEVDepth	#1	-	662.6	×	-	-	0.481	0.383
	#1	8.5	662.6	√	0.221	0.177	0.464	0.327
	#2	17.2	662.6	√	0.210	0.152	0.422	0.289
	#3	11.0	662.6	√	0.218	0.183	0.501	0.332
BEVDepth-GM	#1	11.8	662.6	√	0.221	0.177	0.466	0.348
	#2	24.9	662.6	√	0.210	0.152	0.421	0.319
	#3	12.8	662.6	√	0.218	0.183	0.511	0.355

- (1) These 3D detectors evaluated show a significant performance drop on the ASIO benchmark compared to the offline evaluation. When equipped with camera #1, the 3D detectors BEVFormer, FCOS3D and PGD, which have high computational resource requirements (GFLOPs>1300), suffer a decrease of 20.8%, 27.0% and 33.7% in mAP-A respectively compared to the offline evaluation. For the efficient detectors (frame rate>8) BEVDepth, BEVDet and BEVDet4D, the mAP-A still drops by 1.1%, 1.6% and 6.8%.
- (2) For the evaluation of streaming inputs metrics \mathcal{H} and \mathcal{D} , they can serve as an indicator of the inference speed and accuracy of streaming perception. As the metric \mathcal{H} of streaming inputs increases from 0.210@#2 to 0.221@#1, the inference frame rates of FCOS3D, PGD, BEVDet, BEVDet4D, BEVFormer, PETR, BEVDepth drop by 53.1%, 50.0%, 28.1%, 37.8%, 58.0%, 56.9%, 50.6% respectively. And as the metric \mathcal{D} increases from 0.0152@#2 to 0.0183@#3, the mAP-A of above methods increased by 40.3%, 62.1%, 7.5%, 6.1%, 45.1%, 12.1%, 9.7%. Observations reveal that among all detectors subjected to our testing, there exists an inverse correlation between the \mathcal{H} -value and inference speed when deploying detectors with different combinations of major cameras. Conversely, a positive correlation is observed between the \mathcal{D} -value and online detection accuracy. Therefore, predictive assessments of their real-time inference timeliness and accuracy during actual deployment can be achieved through pre-established \mathcal{H} and \mathcal{D} indices.
- (3) Under different types of major cameras, the magnitude of model performance variation varies widely. As illustrated in Figure 1, a substantial decline in mAP-A is evident for the efficient models BEVDet, BEVDet4D, and PETR when transitioning from high-definition cameras #1 and #3 to the wide-angle camera #2. Specifically, the mAP-A values experienced significant reductions of at least 31.1%, 17.8%, and 31.4%, respectively. In contrast, models with high computational resource requirements, FCOS3D and PGD, exhibited relatively modest performance decrements, amounting to 6.1% and 8.0%, respectively. Notably, the accuracy rankings of these two models were inverted between offline and online testing. This inversion underscores the greater practical significance of ASIO benchmarking for the actual deployment of perception system as compared to offline testing.
- (4) The Backbone switcher can modulate the load on the model's computational resources to compensate for inference delays, thus improving streaming perception. The BEVDepth-GM, equipped with our backbone switcher, demonstrated respective mAP-A improvements of 6.4%, 10.3%, and 6.9% on major cameras #1, #2, and #3. Additional test results presented in Table 7 indicate that FCOS3D, PGD, and BEVFormer achieved mAP-A enhancements of 5.6%, 5.2%, and 9.9%, respectively, on camera #1. It is noteworthy that, owing to its lower model complexity, BEVFormer-GM exhibited more significant improvements (i.e., GFLOPs@BEVFormer is 1322.2,

which is significantly lower than that of FCOS3D (2008.2) and PGD (2223.0)). This underscores the efficacy of switching the backbone in practical deployment scenarios, particularly for simpler models. Furthermore, our observations reveal that the backbone switcher has a more pronounced impact on the AP-A of high-speed objects. For example, on BEVFormer@#1, the AP-A for car and bus increased by 7.3% and 7.9%, respectively, whereas the AP-A for slow-speed objects (pedestrian) saw an increase of 1.7%. This insight can inform future streaming algorithms to concurrently consider major camera selection and the speed differences among different object categories.

Table 7. mAP-A of FCOS3D, PGD, BEVFormer and the corresponding models with backbone switchers. The experiments are conducted under the deployment of major camera #1 while we report AP-A for high-speed categories (e.g., cars, buses) and slow-speed categories (e.g., pedestrians).

Method	mAP-A \uparrow	AVE \downarrow	Car	Bus	Ped.
FCOS3D	0.214	1.309	0.251	0.110	0.301
FCOS3D-GM	0.226(+5.6%)	1.298	0.273	0.132	0.315
PGD	0.212	1.284	0.245	0.102	0.297
PGD-GM	0.223(+5.2%)	1.280	0.269	0.127	0.308
BEVFormer	0.332	0.391	0.381	0.316	0.415
BEVFormer-GM	0.365(+9.9%)	0.388	0.409	0.341	0.422

The aforementioned experimental findings underscore the substantial impact of different types of inputs on the performance of streaming perception. Although efficient detectors such as BEVDet4D, BEVDet, and PETR exhibit excellent performance in offline testing and high-resolution inputs, their performance experiences significant degradation when exposed to wide-field-of-view streaming inputs. In contrast, complex models such as FCOS3D, PGD, and BEVDepth demonstrate more consistent performance across various streaming inputs. Furthermore, by conducting pre-assessments of complexity through \mathcal{H} and \mathcal{D} evaluations for streaming inputs from different major cameras, one can effectively estimate the efficiency and accuracy of streaming perception for these 3D detectors under such input conditions.

4.3. Analysis on Computational Resources Sharing

Typically, the same major camera may be employed for multiple tasks as shared streaming inputs. To analyze the performance fluctuations induced by shared computational resources, we evaluated the 3D detectors (BEVFormer and BEVDepth) on a GPU (RTX4090) concurrently processing N classification tasks based on ResNet18. As illustrated in Figure 7, as the number of classification tasks increases, the performance of BEVFormer and BEVDepth declines due to the reduction in computational resources allocated to 3D detection tasks. Specifically, as the number of classification tasks increases from 0 to 10, the mAP-S of BEVFormer and BEVDepth decrease by 49.5% and 20.2%, respectively. It is noteworthy that the proposed backbone switcher consistently enhances streaming performance under computational sharing conditions. When executing 10 classification tasks, the mAP-A of BEVDepth-GM and BEVFormer-GM increased by 10.9% and 15.9%, respectively.

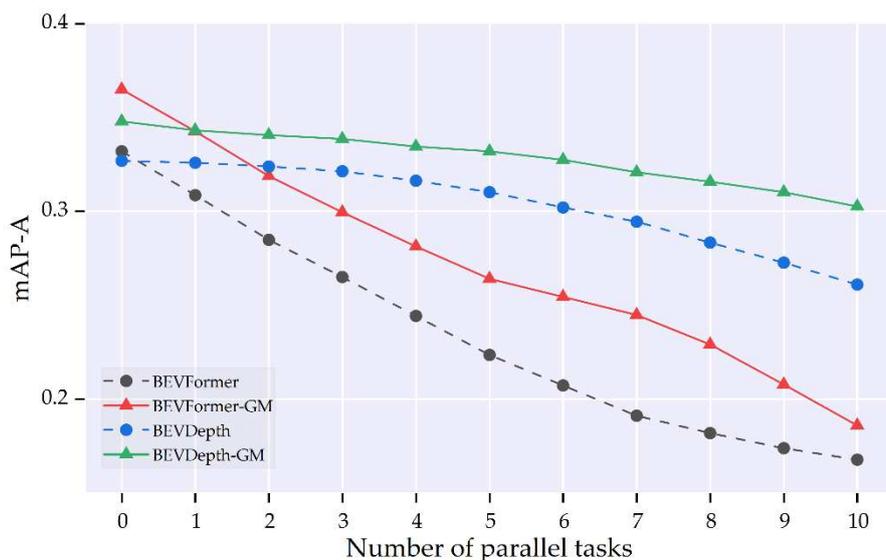


Figure 7. Comparison of the streaming performance of BEVFormer, BEVFormer-GM, BEVDepth, and BEVDepth-GM on major camera #1 under computational resources sharing. The x-axis denotes the quantity of parallel ResNet18-based classification tasks.

5. Discussion and Conclusions

This paper introduces the ASIO benchmark for evaluating the online performance of vision-centric autonomous driving perception systems under various streaming inputs. Specifically, we establish an evaluation dataset based on the CARLA Leaderboard, serving dual purposes: estimating the information complexity of streaming inputs in advance and validating camera-based streaming 3D detection. The evaluation metrics encompass two components—an information complexity assessment metric involving a fractional-dimensional two-dimensional entropy specifically tailored to input information from different major cameras, and a performance evaluation metric based on ground truth distinct from offline methods. Additionally, we propose the ASIO switcher based on the real-time input's information complexity to address abrupt changes in input information for 3D detectors, consistently achieving superior streaming performance across three major cameras. Leveraging the ASIO benchmark, we investigate the online performance of seven representative 3D detectors under different streaming inputs. The experimental results demonstrate that the information complexity of streaming inputs can be utilized to predict the practical deployment online performance of 3D detectors. Furthermore, considerations of the model's parallel computational budget and the selection of backbones based on varying information complexities should be incorporated into the design considerations for practical deployment. While the proposed ASIO benchmark represents a significant stride towards practical vision-centric perception in autonomous driving, several limitations warrant further investigation in future research: (1) the establishment of more comprehensive and enriched datasets is needed to adequately address performance testing of streaming perception from input to algorithm and computational platforms; (2) in real-world deployment, an extremely diverse sensor configuration is adopted, encompassing multi-camera setups, infrared sensors, and even event cameras, necessitating the development of a more generalized and unified description for such configurations; (3) the assessment of real-time inputs and corresponding strategies for 3D detectors merit further research; (4) algorithms geared towards multitasking and end-to-end approaches should encompass a broader spectrum of autonomous driving tasks, such as depth estimation and dynamic tracking, requiring inclusion in the computation of evaluation metrics.

Author Contributions: Conceptualization, T.J. and W.D.; Formal analysis, M.Y.; Investigation, T.J.; Methodology, T.J.; Resources, M.Y.; Supervision, W.D. and M.Y.; Validation, T.J., H.Z. and P.D.; Visualization,

P.D and H.Z.; Writing—original draft, T.J.; Writing—review and editing, T.J. and W.D. All authors have read and agreed to the published version of the manuscript.

Funding: This work was supported by the National Natural Science Foundation of China under Grant/Award Number 51775451; the Natural Science Foundation of Sichuan Province under Grant/Award Numbers 2023NSFSC0395 and 2022NSFSC1892; and the Sichuan Science and Technology Program under Grant/Award Number 2022ZHCG0061.

Data Availability Statement: Data sharing not applicable.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Huang, J.; Huang, G.; Zhu, Z.; Ye, Y.; Du, D. Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. arXiv preprint arXiv:2112.11790 2021.
2. Huang, J.; Huang, G. Bevdet4d: Exploit temporal cues in multi-camera 3d object detection. arXiv preprint arXiv:2203.17054 2022.
3. Li, Y.; Bao, H.; Ge, Z.; Yang, J.; Sun, J.; Li, Z. Bevstereo: Enhancing depth estimation in multi-view 3d object detection with temporal stereo. In Proceedings of the Proceedings of the AAAI Conference on Artificial Intelligence, 2023; pp. 1486-1494.
4. Li, Y.; Ge, Z.; Yu, G.; Yang, J.; Wang, Z.; Shi, Y.; Sun, J.; Li, Z. Bevdepth: Acquisition of reliable depth for multi-view 3d object detection. In Proceedings of the Proceedings of the AAAI Conference on Artificial Intelligence, 2023; pp. 1477-1485.
5. Li, Z.; Wang, W.; Li, H.; Xie, E.; Sima, C.; Lu, T.; Qiao, Y.; Dai, J. Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers. In Proceedings of the European conference on computer vision, 2022; pp. 1-18.
6. Liu, Y.; Wang, T.; Zhang, X.; Sun, J. Petr: Position embedding transformation for multi-view 3d object detection. In Proceedings of the European Conference on Computer Vision, 2022; pp. 531-548.
7. Liu, Y.; Yan, J.; Jia, F.; Li, S.; Gao, A.; Wang, T.; Zhang, X. Petr2: A unified framework for 3d perception from multi-camera images. In Proceedings of the Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023; pp. 3262-3272.
8. Jiang, Y.; Zhang, L.; Miao, Z.; Zhu, X.; Gao, J.; Hu, W.; Jiang, Y.-G. Polarformer: Multi-camera 3d object detection with polar transformer. In Proceedings of the Proceedings of the AAAI Conference on Artificial Intelligence, 2023; pp. 1042-1050.
9. Li, Y.; Chen, Y.; Qi, X.; Li, Z.; Sun, J.; Jia, J. Unifying voxel-based representation with transformer for 3d object detection. *Advances in Neural Information Processing Systems* 2022, 35, 18442-18455.
10. Park, J.; Xu, C.; Yang, S.; Keutzer, K.; Kitani, K.; Tomizuka, M.; Zhan, W. Time will tell: New outlooks and a baseline for temporal multi-view 3d object detection. arXiv preprint arXiv:2210.02443 2022.
11. Li, Q.; Wang, Y.; Wang, Y.; Zhao, H. Hdmapnet: An online hd map construction and evaluation framework. In Proceedings of the 2022 International Conference on Robotics and Automation (ICRA), 2022; pp. 4628-4634.
12. Pan, B.; Sun, J.; Leung, H.Y.T.; Andonian, A.; Zhou, B. Cross-view semantic segmentation for sensing surroundings. *IEEE Robotics and Automation Letters* 2020, 5, 4867-4873.
13. Peng, L.; Chen, Z.; Fu, Z.; Liang, P.; Cheng, E. BEVSegFormer: Bird's Eye View Semantic Segmentation From Arbitrary Camera Rigs. In Proceedings of the Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2023; pp. 5935-5943.
14. Liu, Y.; Yuan, T.; Wang, Y.; Wang, Y.; Zhao, H. Vectormapnet: End-to-end vectorized hd map learning. In Proceedings of the International Conference on Machine Learning, 2023; pp. 22352-22369.
15. Liao, B.; Chen, S.; Wang, X.; Cheng, T.; Zhang, Q.; Liu, W.; Huang, C. Maptr: Structured modeling and learning for online vectorized hd map construction. arXiv preprint arXiv:2208.14437 2022.
16. Akan, A.K.; Güneş, F. Stretchbev: Stretching future instance prediction spatially and temporally. In Proceedings of the European Conference on Computer Vision, 2022; pp. 444-460.
17. Hu, A.; Murez, Z.; Mohan, N.; Dudas, S.; Hawke, J.; Badrinarayanan, V.; Cipolla, R.; Kendall, A. Fiery: Future instance prediction in bird's-eye view from surround monocular cameras. In Proceedings of the Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021; pp. 15273-15282.

18. Kendall, A.; Hawke, J.; Janz, D.; Mazur, P.; Reda, D.; Allen, J.-M.; Lam, V.-D.; Bewley, A.; Shah, A. Learning to drive in a day. In Proceedings of the 2019 International Conference on Robotics and Automation (ICRA), 2019; pp. 8248-8254.
19. Arshad, S.; Suaaleh, M.; Kim, D.; Nam, D.V.; Kim, G.-W. Clothoid: an integrated hierarchical framework for autonomous driving in a dynamic urban environment. *Sensors* 2020, 20, 5053.
20. Zhu, Z.; Zhao, H. Learning Autonomous Control Policy for Intersection Navigation With Pedestrian Interaction. *IEEE Transactions on Intelligent Vehicles* 2023.
21. Lang, A.H.; Vora, S.; Caesar, H.; Zhou, L.; Yang, J.; Beijbom, O. Pointpillars: Fast encoders for object detection from point clouds. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019; pp. 12697-12705.
22. Yan, Y.; Mao, Y.; Li, B. Second: Sparsely embedded convolutional detection. *Sensors* 2018, 18, 3337.
23. Yin, T.; Zhou, X.; Krahenbuhl, P. Center-based 3d object detection and tracking. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2021; pp. 11784-11793.
24. Caesar, H.; Bankiti, V.; Lang, A.H.; Vora, S.; Liong, V.E.; Xu, Q.; Krishnan, A.; Pan, Y.; Baldan, G.; Beijbom, O. nuscenes: A multimodal dataset for autonomous driving. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020; pp. 11621-11631.
25. Cordts, M.; Omran, M.; Ramos, S.; Rehfeld, T.; Enzweiler, M.; Benenson, R.; Franke, U.; Roth, S.; Schiele, B. The cityscapes dataset for semantic urban scene understanding. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2016; pp. 3213-3223.
26. Geiger, A.; Lenz, P.; Urtasun, R. Are we ready for autonomous driving? the kitti vision benchmark suite. In Proceedings of the 2012 IEEE conference on computer vision and pattern recognition, 2012; pp. 3354-3361.
27. Geyer, J.; Kassahun, Y.; Mahmudi, M.; Ricou, X.; Durgesh, R.; Chung, A.S.; Hauswald, L.; Pham, V.H.; Mühlegg, M.; Dorn, S. A2d2: Audi autonomous driving dataset. arXiv preprint arXiv:2004.06320 2020.
28. Huang, X.; Wang, P.; Cheng, X.; Zhou, D.; Geng, Q.; Yang, R. The apollo open dataset for autonomous driving and its application. *IEEE transactions on pattern analysis and machine intelligence* 2019, 42, 2702-2719.
29. Neuhold, G.; Ollmann, T.; Rota Bulò, S.; Kotschieder, P. The mapillary vistas dataset for semantic understanding of street scenes. In Proceedings of the Proceedings of the IEEE international conference on computer vision, 2017; pp. 4990-4999.
30. Scheel, O.; Bergamini, L.; Wolczyk, M.; Osiński, B.; Ondruska, P. Urban driver: Learning to drive from real-world demonstrations using policy gradients. In Proceedings of the Conference on Robot Learning, 2022; pp. 718-728.
31. Sun, P.; Kretschmar, H.; Dotiwalla, X.; Chouard, A.; Patnaik, V.; Tsui, P.; Guo, J.; Zhou, Y.; Chai, Y.; Caine, B. Scalability in perception for autonomous driving: Waymo open dataset. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020; pp. 2446-2454.
32. Wilson, B.; Qi, W.; Agarwal, T.; Lambert, J.; Singh, J.; Khandelwal, S.; Pan, B.; Kumar, R.; Hartnett, A.; Pontes, J.K. Argoverse 2: Next generation datasets for self-driving perception and forecasting. arXiv preprint arXiv:2301.00493 2023.
33. Yu, F.; Chen, H.; Wang, X.; Xian, W.; Chen, Y.; Liu, F.; Madhavan, V.; Darrell, T. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In Proceedings of the Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020; pp. 2636-2645.
34. Zhang, S.; Benenson, R.; Schiele, B. Citypersons: A diverse dataset for pedestrian detection. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2017; pp. 3213-3221.
35. Wang, X.; Zhu, Z.; Zhang, Y.; Huang, G.; Ye, Y.; Xu, W.; Chen, Z.; Wang, X. Are We Ready for Vision-Centric Driving Streaming Perception? The ASAP Benchmark. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023; pp. 9600-9610.
36. Carla autonomous driving leaderboard (accessed November 2021). <https://leaderboard.carla.org/leaderboard/>, 2021
37. Li, M.; Wang, Y.-X.; Ramanan, D. Towards streaming perception. In Proceedings of the Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16, 2020; pp. 473-488.
38. Kalman, R.E. A new approach to linear filtering and prediction problems. 1960.

39. Ghosh, A.; Nambi, A.; Singh, A.; Yvs, H.; Ganu, T. Adaptive streaming perception using deep reinforcement learning. arXiv preprint arXiv:2106.05665 2021.
40. Han, W.; Zhang, Z.; Caine, B.; Yang, B.; Sprunk, C.; Alsharif, O.; Ngiam, J.; Vasudevan, V.; Shlens, J.; Chen, Z. Streaming object detection for 3-d point clouds. In Proceedings of the European Conference on Computer Vision, 2020; pp. 423-441.
41. Peng, C.-K.; Buldyrev, S.V.; Havlin, S.; Simons, M.; Stanley, H.E.; Goldberger, A.L. Mosaic organization of DNA nucleotides. *Physical review e* 1994, 49, 1685
42. Warfield, J.N. Societal systems planning, policy and complexity. *Cybernetics and System* 1978, 8, 113-115.
43. Rényi, A. On measures of entropy and information. In Proceedings of the Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics, 1961; pp. 547-562.
44. Tsallis, C. Possible generalization of Boltzmann-Gibbs statistics. *Journal of statistical physics* 1988, 52, 479-487.
45. Pincus, S.M. Approximate entropy as a measure of system complexity. *Proceedings of the National Academy of Sciences* 1991, 88, 2297-2301.
46. Pincus, S.M.; Goldberger, A.L. Physiological time-series analysis: what does regularity quantify? *American Journal of Physiology-Heart and Circulatory Physiology* 1994, 266, H1643-H1656.
47. Richman, J.S.; Moorman, J.R. Physiological time-series analysis using approximate entropy and sample entropy. *American journal of physiology-heart and circulatory physiology* 2000, 278, H2039-H2049.
48. Ubriaco, M.R. Entropies based on fractional calculus. *Physics Letters A* 2009, 373, 2516-2519.
49. Machado, J.T. Fractional order generalized information. *Entropy* 2014, 16, 2350-2361.
50. Wang, T.; Zhu, X.; Pang, J.; Lin, D. Fcos3d: Fully convolutional one-stage monocular 3d object detection. In Proceedings of the Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021; pp. 913-922.
51. Wang, Y.; Guizilini, V.C.; Zhang, T.; Wang, Y.; Zhao, H.; Solomon, J. Detr3d: 3d object detection from multi-view images via 3d-to-2d queries. In Proceedings of the Conference on Robot Learning, 2022; pp. 180-191.
52. Li, E.; Wang, S.; Li, C.; Li, D.; Wu, X.; Hao, Q. Sustech points: A portable 3d point cloud interactive annotation platform system. In Proceedings of the 2020 IEEE Intelligent Vehicles Symposium (IV), 2020; pp. 1108-1115.
53. Qian, X.; Liu, C.; Qi, X.; Tan, S.-C.; Lam, E.; Wong, N. Context-Aware Transformer for 3D Point Cloud Automatic Annotation. arXiv preprint arXiv:2303.14893 2023.
54. Xiong, H.; Shang, P.; Zhang, Y. Fractional cumulative residual entropy. *Communications in Nonlinear Science and Numerical Simulation* 2019, 78, 104879.
55. Liu, S.; Chunwu, Y.; Dazhi, C. Weapon equipment management cost prediction based on forgetting factor recursive GM (1, 1) model. *Grey Systems: Theory and Application* 2020, 10, 38-45.
56. Wang, T.; Xinge, Z.; Pang, J.; Lin, D. Probabilistic and geometric depth: Detecting objects in perspective. In Proceedings of the Conference on Robot Learning, 2022; pp. 1475-1485.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.