

Generating Realistic and Diverse Tests for LiDAR-Based Perception Systems

Garrett Christian
University of Virginia
Charlottesville, Virginia, USA
rda2tc@virginia.edu

Trey Woodlief*
University of Virginia
Charlottesville, Virginia, USA
adw8dm@virginia.edu

Sebastian Elbaum
University of Virginia
Charlottesville, Virginia, USA
selbaum@virginia.edu

Abstract—Autonomous systems rely on a perception component to interpret their surroundings, and when misinterpretations occur, they can and have led to serious and fatal system-level failures. Yet, existing methods for testing perception software remain limited in both their capacity to efficiently generate test data that translates to real-world performance and in their diversity to capture the long tail of rare but safety-critical scenarios. These limitations are particularly evident for perception systems based on LiDAR sensors, which have emerged as a crucial component in modern autonomous systems due to their ability to provide a 3D scan of the world and operate in all lighting conditions. To address these limitations, we introduce a novel approach for testing LiDAR-based perception systems by leveraging existing real-world data as a basis to generate realistic and diverse test cases through mutations that preserve realism invariants while generating inputs rarely found in existing data sets, and automatically crafting oracles that identify potentially safety-critical issues in perception performance. We implemented our approach to assess its ability to identify perception failures, generating over 50,000 test inputs for five state-of-the-art LiDAR-based perception systems. We found that it efficiently generated test cases that yield errors in perception that could result in real consequences if these systems were deployed and does so at a low rate of false positives.

Index Terms—Software Testing and Validation, Machine Learning

I. INTRODUCTION

As autonomous system (AS) research and development accelerates, so too must our efforts to rigorously test them to reduce field failures. A key component underpinning all aspects of ASs is the perception software that transforms sensor inputs into an interpretation of the world [18]. Since all actions of the AS are predicated on its world interpretation, perception failures can lead to safety-critical system-level failures. Field failures in AS perception components have resulted in multiple fatalities in recent years [4], [5], [44]. For example, in 2018 an Uber autonomous vehicle going 35 mph at night collided with a woman crossing the street with her bike [16], [19], [32]. The perception component consuming a LiDAR point cloud (PC), a vector containing a set of locations in 3D denoting how far away an object is from the sensor (as exemplified in Fig. 1a), produced an interpretation of the bike that oscillated between vehicle, “other”, and bike in the 10 seconds leading up to the crash [4]. Only 1.2 seconds before

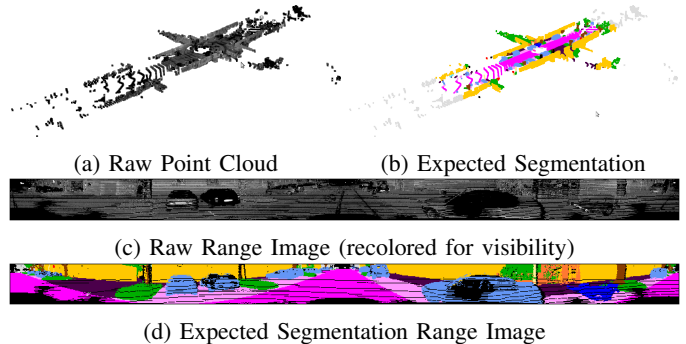


Fig. 1: Example LiDAR Point Cloud (PC) used in perception (best viewed on a screen). Fig. 1a shows a PC in perspective. Figs. 1a & 1c are colored with the reflected LiDAR beam intensities. Figs. 1b & 1d are colored with the human-annotated semantic class. Figs. 1c & 1d show a “range image” where the 360° view is projected onto a cylinder and unrolled.

the collision did the component produce a consistently correct interpretation, which was too late to alert the human driver or brake, resulting in a fatal collision.

Exposing these failures, however, can be particularly challenging when the input is provided by LiDAR. While LiDAR is becoming more prevalent in AS companies (e.g., Volvo [31], Motional [3], Baidu [54], Ford [1], Zoox [37], and Waymo [36]), techniques to generate test cases that adequately reflect the safety-critical scenarios ASs face in the field are limited along two critical dimensions: realism and diversity.

First, although real-world LiDAR PCs can be collected while driving, generating their expected outputs through manual labeling can be prohibitively expensive as PCs are difficult to label, requiring significant human effort along with specialized training and software, limiting availability to a few data sets [2], [6], [9], [23], [35], [38], [42]. As an example of the difficulty in labeling, Fig. 1 shows a sample PC presented as a 3D scene in Fig. 1a and as a 2D projection in Fig. 1c with their respective expected semantic labels in Figs. 1b and 1d. While large objects such as vehicles may be recognizable by their shape, discerning labels for other points is challenging due to the disconnect between PCs and how humans perceive the world—we have the analog of built-in cameras, but not

*Corresponding Author

built-in LiDAR. Besides labeling challenges, these data sets are inherently *limited in diversity* as the collected data is unlikely to include dangerous edge cases such as a person crossing directly in front of the vehicle. This rarity is a fundamental limit of data collection efforts as it is unfeasible, uneconomical, and unethical to collect the safety-critical edge cases necessary for sufficient diversity in testing.

Second, although simulation-based techniques can serve to generate synthetic LiDAR inputs and use the simulator state to compute the expected outputs [13], [21], [47], [56], the simulation-reality domain gap [27] makes it difficult to judge whether a system’s performance is due to the inherent synthesised nature of the data or due to component failures [29], [47]. When a failure is uncovered in simulation, is it a failure that could not occur in the real world, i.e. a false positive? In spite of recent advances in LiDAR simulation techniques to mitigate this challenge by supplementing simulators with real-world PCs [30], [48], the simulation-reality gap remains a significant obstacle. Although simulation allows for improved test diversity, it is *limited by the level of realism*.

Finally, LiDAR-specific data generation techniques have not focused on the specific need to generate realistic data from the long tail of possible scenes in order to increase diversity. The two closest lines of work are as follows. First are techniques for generating adverse weather effects [24], [28] and whole-PC affine transformations [22] that, while realistic, provide *limited increases in diversity* as they do not change expected output at the semantic level, i.e. the semantics of each point is the same even if the point is moved or obscured. Second are techniques for adding objects to a scene [15], [26], [53] that have the potential to increase diversity by altering scene semantics, but *do not ensure realism*. Section V-B contains a full discussion of related approaches.

Recent work on testing image-based perception systems proposed using existing real-world images and performing mutations at the semantic level to alter both the input and expected output [50]. Similar to our aim, this work focuses on realism and diversity, generating realistic tests that are often lacking in existing data sets. While this set the stage for sensor mutation as a test generation technique regardless of the sensor used, its image-centric focus led to three crucial issues limiting its direct applicability to other sensor domains: over-dependence on the image domain to define mutation operators, reliance on human perception for false positive detection, and inability to adequately judge the performance of SUTs.

Building on this work on semantic mutation testing for images, we perform the non-trivial extension to generate test cases for LiDAR-based perception systems, aiming to generate test cases that explore the long tail of possible scenes, e.g. the case of a woman walking with her bike in front of the AS. Different from images, semantic mutations for LiDAR require a reimagining of the mutation space, with special considerations for the physics of the LiDAR sensor and the 3D nature of the input domain, and specialized techniques for judging realism and evaluating the performance of the system under test (SUT) as even small failures such as misclassifying

a pedestrian that is 0.05% of the total PC can be safety-critical.

Our approach pursues realism by starting with training data consisting of real-world PCs that have been semantically labeled, tapping into this real data to build resources that can be used in mutation operations, and by designing domain-specific mutations that respect real-world invariants. Our approach pursues diversity by developing mutations that can complement the available real-world data sets to produce diverse test inputs that could be found in real life situations. Then, by using metamorphic relationships between the original and mutated data, we can find potentially dangerous failures in perception that are not covered by the original data set. Further, motivated by the inherent difficulty for humans to perform the LiDAR perception task, this work proposes an automated mechanism for evaluating the realism of the mutations; by examining the performance of multiple SUTs, we can automatically identify potential false positives from unrealistic inputs. Our contributions are as follows:

- A new approach for testing LiDAR-based perception systems, based on semantic mutations applied to labeled real-world data.
- An implementation targeting LiDAR-based perception of autonomous vehicles using the SemanticKITTI dataset [2], publicly releasing the code¹.
- Experiments using our approach to test five state-of-the-art LiDAR-based perception systems, revealing errors in perception that could result in real consequences if these systems were deployed.

II. APPROACH

Our approach takes as input a data set of LiDAR PCs collected from the real world and labeled with ground truth semantics, and performs realistic semantic mutations to generate a diverse set of novel test cases that meaningfully exercise LiDAR perception systems. We respect realism by starting with real-world tests, collecting entities from the real world as resources, and crafting realism invariants that must hold during mutation. We increase diversity by crafting mutations that can render hard to find cases, such as a pedestrian in the roadway. Our approach occurs in three stages which we discuss in detail in this section: Resource Collection, Mutation Generation, and Performance Evaluation.

A. Preliminaries

LiDAR, or **L**ight **D**etection **A**nd **R**anging, technology measures the round-trip time-of-flight of a beam of light between the sensor and an object to determine the distance to the object. Many LiDAR systems used in ASs leverage a rotating sensor that allows the system to take discrete measurements at regular intervals to obtain the distance to many different objects within a scene. Taken together, these measurements form a *point cloud* (PC), a representation of the sensor reading that uses the sensed distances and calculated angles to create a set of points in 3D space that depict the locations of all sensed objects.

¹<https://github.com/less-lab-uva/semLidarFuzz>

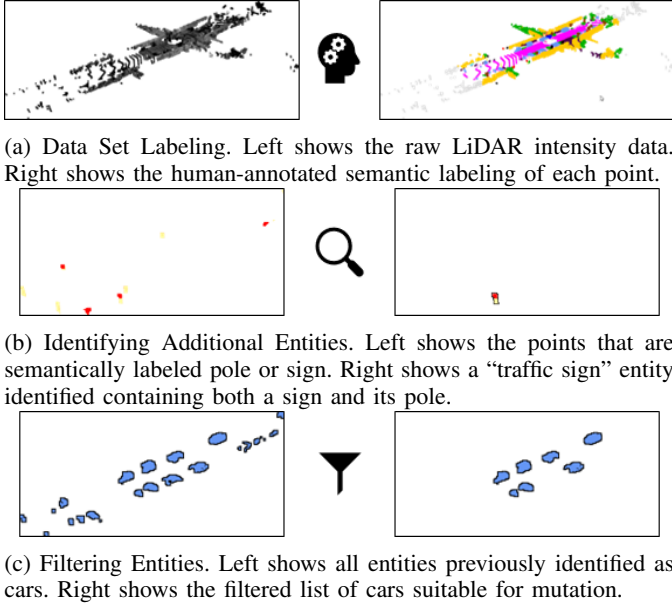


Fig. 2: Stage 1: Resource Collection

In many ASs, the LiDAR serves as an input to the perception software, leveraging the PC to perceive the semantics of the scene. We consider the semantic segmentation task in which a perception system assigns each point in a PC a semantic class identifying what the point represents, e.g. road, person, car, etc. A test case, $t = (PC, ExpSem)$, for a LiDAR-based perception system consists of an input PC, PC , an expected output semantic labeling, $ExpSem$, and an evaluation function, $eval$. The label produced by the SUT, $SUT(PC)$, is evaluated against the expected output using the evaluation function, i.e. $eval(ExpSem, SUT(PC))$ to determine if the SUT passes the test.

As many AS perception systems use machine-learned components, test cases can also be used to train the system. This has given rise to several data sets that consist of LiDAR-perception tests containing PC s gathered from the real world with accompanying $ExpSem$ from human annotation for vehicle-centric scenes (e.g., [2], [6], [9], [35], [38], [42]); additionally, some data sets (e.g., [2], [6], [35], [42]) also provide an *entity label* for each point, describing what logical group each point belongs to, e.g. a set of points corresponding to a particular vehicle. Due to the expense of collecting and annotating these data sets, they are limited in scope. Starting from this one-time cost to obtain, we can leverage the existing data set to produce novel and relevant test cases—in our study we more than tripled the available tests using our technique.

B. Resource Collection

The first stage of our approach leverages the provided data set to generate a resource set for the later mutation operations. The resource set contains a list of PCs suitable for mutation, the semantic labeling of each PC, and the list of entities, E , appearing in each PC. Resource Collection occurs in two phases, illustrated in Fig. 2.

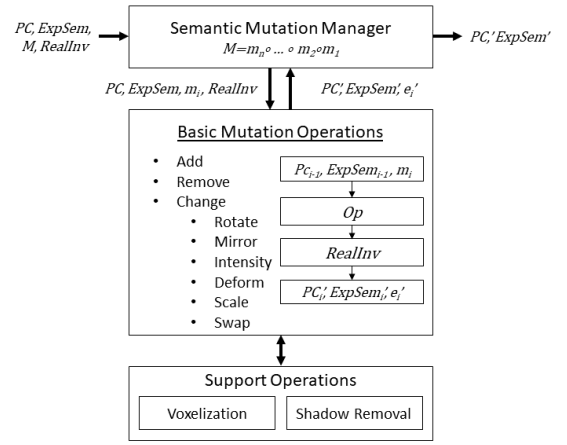


Fig. 3: Stage 2: Input Generation

First, the Entity Identification phase (Fig. 2b) attempts to find additional entities that were not initially recorded in the data set. This is useful due to the different entity identification paradigms used to generate existing real-world data sets in order to provide a rich and consistent entity set. For example, the SemanticKITTI data set [2] does not include traffic signs in its entity lists; however, it does label points with a traffic sign class. Thus, a traffic sign entity can be reconstructed by identifying a collection of neighboring labeled traffic sign points. For each of the provided $(PC, ExpSem, E_{orig})$ in the data set, the semantic information of $ExpSem$ is used along with the spatial relationships in PC and a set of parameterized invariants, further described in Section III-A, to identify additional entities that were not provided with the initial entity list, E_{add} , resulting in an expanded $E_{total} = E_{orig} \cup E_{add}$.

Second, each E_{total} is filtered (Fig. 2c) to identify entities suitable for mutation in generating test cases: $E_{filtered} = Filter(E_{total})$. Entities that are too small, too far from the sensor, or partially obscured are removed. This process is parameterized and helps ensure that the later stages of test case generation provide realistic test inputs. Removing entities that are too small or too far from the sensor helps to ensure that the entities have sufficient definition and are not affected by sensor noise. Removing entities that are partially obscured ensures that the points that form the entity are complete, e.g. the car is not missing its hood because it was obscured by another car. Once complete, the Resource Collection stage outputs a resource set consisting of a set of $(PC, ExpSem, E_{filtered})$ that are suitable for mutation. This process only needs to be completed once for a given data set.

C. Mutation Generation

The second stage of the technique focuses on generating a realistic mutation, Mut , to transform an existing PC into a novel input with its expected semantics. Mut takes as input a PC, PC , its expected semantics, $ExpSem$, and a mutation parameter set, M , along with a set of parameterized realism invariants, $RealInv$. Let $t = (PC, ExpSem)$ be

an initial test case that is mutated to form a new test case $t' = (PC', ExpSem') = Mut(PC, ExpSem, M, RealInv)$. For brevity, in the context of a single mutation, we write $t' = Mut(t)$, etc. M consists of what type of mutation should be performed along with any mutation-specific parameters, such as what entities from the resource set should be used. $RealInv$ is a parameterized set of realism invariants that must hold for the final t' to be realistic. During the mutation process, if an invariant is violated, the mutation will abort and a new test will not be generated; this means that the provided $PC, ExpSem, M$ were incompatible with $RealInv$. We adopt the definition of realism used in prior work on semantic mutation for images that a test case is realistic if $PC, ExpSem$ could come from a possible configuration of the real world, including the long tail of improbable configurations [50].

This stage is broken into several components as outlined in Fig. 3. First, the Semantic Mutation Manager breaks the mutation M into a composition of basic mutation operations at the entity level, each with their own parameters: $M = m_n \circ \dots \circ m_2 \circ m_1$. The manager then applies and composes each of these basic operations in turn to generate the final $PC', ExpSem'$ pair. Each of the basic operations falls into one of three categories: add, remove, or change, and produces intermediate results for the Semantic Mutation Manager by either updating PC and $ExpSem$, or by altering an entity, e' , to be used in a later operation, e.g. Mirror first produces a mirrored e' that Add then inserts into PC and $ExpSem$. The basic mutation operations further rely on a set of support operations that are used to generate the intermediate PC'_i and $ExpSem'_i$; each basic operation also individually validates that the realism invariants are respected.

1) *Basic Mutation Operations*: Each of the basic mutation operations focus on a single entity and emulate a physical change that could occur in the real world to ensure realism. Each operation has specific invariants that we have recognized as key concerns for that operation which are briefly introduced here, with further discussion on parameterization in Section III. Each mutation is performed on both PC and $ExpSem$ simultaneously so that the mutated $ExpSem$ remains a valid semantic interpretation for the mutated PC . Table I introduces the eight basic mutation operations we explore. While these provide a rich baseline, the “change” category could be further expanded based on desired testing directions as each new operation provides a realistic mechanism to further expand the diversity of the generated tests in a new dimension.

2) *Support Operations*: There are two support operations used by multiple basic mutation operations.

Voxelization: Many invariants, e.g. those of the Rotate operation, rely on checking that an entity’s location is allowable under a given constraint, which can be expensive to validate. Voxelization is a technique rooted in computer graphics [11], [14] that discretizes a 3D representation into an occupancy grid, allowing for efficient computation of functions such as entity intersection or proximity in invariant verification.

Shadow Removal: LiDAR sensors cast rays of light out at regular intervals and use the reflected beam to detect objects.

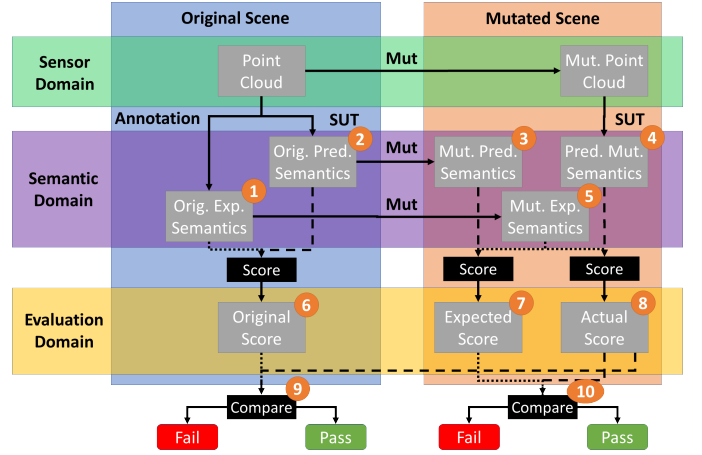


Fig. 4: Stage 3: Performance Evaluation

This means that each object detected by the LiDAR casts a “shadow” as it occludes all objects behind it. When a LiDAR PC is mutated, for example with the Add operation, it is possible for altered points to now occlude existing points, violating this invariant. To restore this property and ensure the invariants are respected, we use a common shadow removal subroutine. The shadow is removed by taking the convex hull of the entity and projecting a mask from the LiDAR sensor to the edges of the entity and then removing any points that fall within that polygon, restoring the shadow.

At the end of this stage, our approach has used Mut to generate a novel test case $t' = (PC', ExpSem')$ that respects the realism invariants provided. At the next stage, the approach will use this information to build the oracle, $eval$, to evaluate the performance of an SUT on this new test.



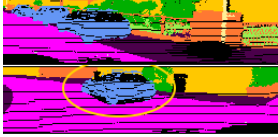
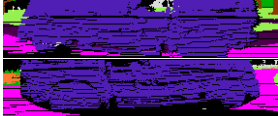
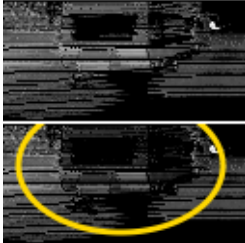
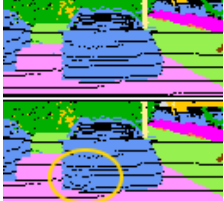
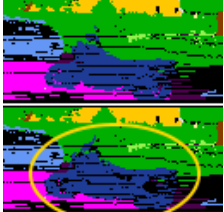
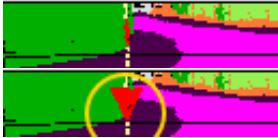
Each mutation faces unique technical challenges in its design and implementation in order to craft and enforce appropriate realism invariants while allowing for efficient production of useful and novel PCs. Table I provides insight about the motivation for utility and notable invariants, while Section III-B provides more information about the technical challenges; additional discussion on both are available in the supplement in the repository.

D. Performance Evaluation

In order to evaluate the performance of the SUTs on the mutated tests, we employ two techniques that allow for more robust and precise performance metrics compared to prior approaches. Fig. 4 outlines two different paradigms for performance evaluation, visualizing the difference between our approach and prior approaches [22], [30], [50]. On the left and right are the original scene and mutated scene respectively; from top to bottom we follow these scenes as they are captured in the sensor domain, perceived in the semantic domain, and then SUT performance is converted to a score in the evaluation domain. The numbered circles correspond to the numbered labels over equation terms later in this section for reference.

In prior literature [22], [30], [50], the score of the SUT, (8), on a mutated test case is based on the intersection over union

TABLE I: Basic Mutation Operations

Name	Category	Motivation	Parameters	Invariants	Example (Before, After)
Add	Add	A crucial limitation of available real-world LiDAR test data is the large imbalance in the semantic classes; many of the most safety-critical classes, e.g. pedestrians, cyclists, etc. are scarce and can be increased through this mutation.	<ul style="list-style-type: none"> Entity to add 	Entity must appear on the ground, and not be occluded by or intersect any other entities. Entity must be placed at the same location in the new <i>PC</i> as it was in its original <i>PC</i> .	
Remove	Remove	SUTs may use the presence of certain entities as context clues when labeling other aspects of the scene, so removing an entity may change the SUTs performance. Once the entity is removed, suitable background points must be added to the PC to fill the space where the entity used to exist.	<ul style="list-style-type: none"> Entity to remove 	The entity to remove must not obscure another entity. Filled in points must be of an appropriate semantic class for the location removed.	
Rotate	Change	Due to the rotational nature of the LiDAR sensor, the location of the entity can rotate around the sensor while maintaining the correct perspective. This movement enables an entity to appear in a context that it did not appear in the data set, e.g. a pedestrian in the roadway.	<ul style="list-style-type: none"> Entity to rotate Angle of new location 	Once in the new position, the entity must not be occluded by scenery.	
Mirror	Change	Common entities in LiDAR scenes can be viewed from a mirrored orientation, e.g. the reverse side of a car. Mirroring increases diversity and rare scenes, such as a vehicle facing the wrong direction on the street. Although likely not present in the training data, it is imperative that a system respond appropriately to this situation.	<ul style="list-style-type: none"> Entity to mirror 	No additional invariants; the mirrored entity still subject to Add invariants.	
Adjust Intensity	Change	LiDAR sensors also detect the intensity of the reflected beam from a surface which can be affected by the material, reflectivity, etc. of the surface. For example, different vehicles naturally have different intensity profiles based on their different paint characteristics and thus this operation is analogous to repainting a car.	<ul style="list-style-type: none"> Entity to adjust intensity Amount to adjust intensity 	Intensity is not adjusted in certain areas, i.e. areas that could not be repainted such as a license plate.	
Deform	Change	Although most vehicles of a given make, model, and year share the same shape, vehicles may be involved in collisions or otherwise accumulate deformations to their exterior; similarly, traffic signs share a common reference shape, but may wear or deform over time. To emulate this, we mutate the entity by deforming a portion of the surface, resembling a dent.	<ul style="list-style-type: none"> Entity to deform Location to deform Size of deformation 	Entity must be deformable, e.g. a vehicle but not a person, and the deformation must be of appropriate size for the entity.	
Adjust Scale	Change	The scale of a class of entity can change over time and across locations. For example, the SemanticKITTI data set [2] was collected in 2012, and in the 10 years since, vehicles have gotten larger. Thus, an entity that is observed in the data set may be similar to an entity that could be encountered in the future, only slightly larger or smaller.	<ul style="list-style-type: none"> Entity to scale Amount to scale 	Entity to scale must contain enough points to have sufficient definition to increase scale; amount to scale must be appropriate for the entity.	
Class-Preserving Swap	Change	The original locations of entities serve as a blueprint for where an entity of a given class could be found in the world, allowing for the replacement of that entity within its class. For example, a stop sign could be replaced with a yield sign. While this does not affect the expected interpretation, this mutation can be used in robustness testing to increase test diversity.	<ul style="list-style-type: none"> Entity to replace Detail of replacement 	The replacement must be of a similar size to the original entity.	

(IoU) metric when comparing the SUT’s PC prediction, ④, with the expected semantics, ⑤. When IoU is calculated over the whole PC, analogous to calculating IoU over the whole image as in prior literature [50], this is a measure of accuracy in terms of the percentage of correct points; later referred to as the “Accuracy metric”. This was also used in past work as the evaluating oracle, ⑨; if the SUT’s IoU score of the mutated test was ϵ below the SUT’s IoU score of the original test, then the SUT failed the test²:

$$\begin{aligned} \text{origScore} &= \text{score}(\overbrace{\text{ExpSem}}^{\textcircled{1}}, \overbrace{\text{SUT(PC)}}^{\textcircled{2}}) \\ \text{mutScore} &= \text{score}(\overbrace{\text{Mut(ExpSem)}}^{\textcircled{5}}, \overbrace{\text{SUT(Mut(PC))}}^{\textcircled{4}}) \\ \epsilon &\overset{\textcircled{9}}{>} \text{origScore} - \text{mutScore} \end{aligned} \quad (1)$$

We start from this baseline and make two key improvements to increase robustness and precision.

1) *Relative Success Metrics*: The robustness of prior approaches is limited by the fact that the mutation changes the test case. Though not previously used as an oracle, when IoU is calculated over only a subset of points, e.g. points belonging to a certain class as in prior LiDAR evaluation [22], [30], this causes the metric to vary as the number of points in the class changes. Thus, in the face of mutations that alter the number of points in a class, the original difference evaluation, ⑨, does not compare like quantities. Consider a scene with no cars that is then mutated to add a car. If the IoU is examining only the “car” class, then the *origScore* for the “car” class is undefined as there is no car. Trying to compare this with the *mutScore* after the car is added does not provide information about the SUT’s performance on the mutated test case.

To address the points changing classes, we adopt a relative success metric that uses the SUT’s original performance on t to calculate an *expected* performance on t' , taking the mutation into account. Under the scheme in Equation 1, the SUT’s performance is evaluated based solely on its separate performance on t (*origScore*) and t' (*mutScore*). However, under the new scheme, we introduce a new term that performs the mutation operation on the SUT’s original prediction for t .

$$\begin{aligned} \text{expScore} &= \text{score}(\overbrace{\text{Mut(ExpSem)}}^{\textcircled{5}}, \overbrace{\text{Mut(SUT(PC))}}^{\textcircled{3}}) \\ \text{mutScore} &= \text{score}(\overbrace{\text{Mut(ExpSem)}}^{\textcircled{5}}, \overbrace{\text{SUT(Mut(PC))}}^{\textcircled{4}}) \\ \epsilon &\overset{\textcircled{10}}{>} \text{expScore} - \text{mutScore} \end{aligned} \quad (2)$$

Fig. 4 visualizes the symmetry of this approach, highlighting how this provides a direct comparison. Returning to the example of the scene with no car being mutated to add a car, the *expScore* has the information about the added car included, allowing the comparison with *mutScore* to provide a meaningful insight into the SUT’s performance.

²The order matters in Equation 1; if $\epsilon > \text{mutScore} - \text{origScore}$, then the mutation improved performance and is not problematic.

2) *Jaccard Index*: Solely relying on an Accuracy metric is ill-suited for LiDAR-based perception systems due to the nature of the entities of interest and possible failure modes; in particular entities “disappearing” or being entirely mislabeled. Consider a mutation that adds a person to a scene. If the SUT mislabels the person, it would represent $\approx 0.05\%$ of the scene’s total points, producing a minuscule drop in accuracy despite its criticality. This implies that a small ϵ is required in order to identify safety-critical issues, yet this does not differentiate what causes the drop in accuracy and could lead to many nuisance failures, e.g. when the SUT mislabels 0.05% of the PC as “terrain” instead of “vegetation”.

This issue is further exacerbated if the mutation changes large portions of the input. Consider a mutation that removes a car. Under Equation 1, if the SUT originally mislabeled the car, it would have a low *origScore*. Then, if the mutation replaced the car with entities the SUT labeled correctly, its *mutScore* would improve by the amount of the car. However, if in the mutated case the SUT also mislabeled a pedestrian, its *mutScore* would still likely be an improvement over its *origScore* since a pedestrian is generally smaller than a car.

As discussed above, prior work used IoU as the scoring metric for evaluating performance. However, the scoring function need not be fixed. The Jaccard Index is a metric commonly used in LiDAR classification evaluation due to its ability to provide a more nuanced metric based on the class label of the affected object [2], [12]. This is particularly useful in the LiDAR context due to the extreme class imbalance in the data; for example, in the SemanticKITTI data set there are roughly 4 orders of magnitude more points labeled “road” than “motorcyclist”. Jaccard Index is calculated by the formula:

$$\frac{1}{C} \sum_{c=1}^C \frac{TP_c}{TP_c + FP_c + FN_c} \quad (3)$$

Here C is the number of classes, the summation iterates over each of the classes c , and TP_c , FP_c , and FN_c are the number of true positive, false positive, and false negative points respectively for the class. Each class contributes equally to the score, allowing for changes in a relatively small class to be reflected in the overall score. Returning to the example of the SUT missing a person, if an SUT mislabels one out of three people in a scene, this is only a $\approx 0.05\%$ drop in overall IoU accuracy; however this would lead to a 33% drop for the “person” class, which out of 20 classes would be a 1.5% drop in the Jaccard Index. This means that for under-represented classes, the Jaccard Index provides a better aggregate score. Together, both accuracy and Jaccard Index provide important metrics of the SUT’s performance on the test, and our approach uses both in tandem.

III. IMPLEMENTATION

We developed an extensible implementation of the approach in Python that automatically generates novel LiDAR test cases through mutation. Each mutation is implemented using Numpy [25], SemanticKITTI’s developer tools [2], and open3D [57].

In this section, due to space, we only describe the implementation of the Resource Set Builder and three representative mutation operations that highlight high-level technical challenges faced by all mutations; a detailed description of all operations is available in the repository.

A. Building the Resource Set

When building the resource set, the filters to identify entities that are too small, too far from the sensor, or occluded are parameterized as follows. First, the entity must contain more than *minPointThreshold* points and the distance to the entity must be less than *distThresh*. Second, to filter occluded objects, the bounding box of the entity is acquired and scaled based on the distance from the sensor by a factor of at most *occlBuffer* to provide a buffer to ensure all occluding objects are identified. The box is also extended to the ground in case another entity is occluding the bottom portion of the object, for example, a car that is only comprised of the roof points. A frustum (bounding cone) is drawn from the sensor to the altered bounding box and if any points that are not semantically classified as unlabeled, outlier, or ground (e.g. road, sidewalk, etc.) are included within this new polygon the entity is discarded. As described in Section II-B, this component also performs Entity Identification to find entities not completely captured in the data set. Since SemanticKITTI does not capture sign entities, the implementation only seeks to identify sign entities. To accomplish this, the points labeled as poles and signs are isolated and grouped into clusters using open3D’s DBSCAN (Density-Based Spatial Clustering of Applications with Noise). If a cluster has both sign points and pole points, has more than *minSignPoints* points, is smaller than *signSizeThreshold*, and meets the other resource set criteria, it is classified as a unique sign entity and added to the entity list. A detailed discussion on parameters and values chosen can be found in the repository.

B. Mutation Technical Challenges

Add Entity. An entity is added into a scene at a given location by adding its points at that location and then removing its shadow. This faces several technical challenges in keeping the added entity in the correct size and pose, verifying that the pose is feasible, and that the surrounding environment is properly occluded. To ensure that the entity appears with the correct size, the distance from the LiDAR sensor is kept consistent. To verify pose and placement when adding an entity, *RealInv* verifies that the entity does not intersect another entity by voxelizing the non-ground points and comparing this to the entity’s new location. The entity is checked for occlusion by drawing a frustum from the entity’s bounding box to the sensor and ensuring that there are fewer than *addOcclusionThreshold* points in the frustum. Finally, *RealInv* verifies that the entity is being placed on the ground in a location that is semantically viable. To do so, the viable ground points are flattened and voxelized and a portion of the entity’s points greater than the *groundThreshold* must be above the ground voxels. Due to the sensor’s point of

view, the ground may not be visible within a short distance of the sensor. To accommodate, if the entity is closer than the *groundDistanceThreshold*, this check is ignored.

Remove Entity. To remove an entity from a scene, the shadow removal process is run in reverse. This faces several technical challenges in ensuring that the points that previously fell on the object are now part of the background behind the object in a continuous and realistic fashion, including its semantic label. To solve this, first the entity’s shadow is generated, showing the area that must be filled to replace the entity. The shadow area is then split in two halves vertically and each half is swept rotationally outward to identify an equal sized portion adjacent to the entity that could be used to back-fill the vacant area. These points are copied into the area, preserving the distance of the points and semantics surrounding the missing area. During this process, *RealInv* verifies that the entity to be removed does not partially obscure another entity and that the points that will fill in the empty area have appropriate semantics. To avoid obscuring another entity, there must be fewer than *removePointsAboveThresh* points above the entity. To ensure appropriate semantics of the filled area, the filled points must not contain points labeled in the *invalidReplacementClasses* set. These classes are chosen since they are likely to result in an unrealistic replacement; e.g., if there are building points in the replacement, these points will be disconnected from the rest of the building.

Rotate Entity. Used in conjunction with the Add Entity mutation, this mutation is used to identify a location within a scene where an entity could be placed and align the entity with that location, while *RealInv* ensures that it is not occluded by scenery. This faces several technical challenges in ensuring that the location slated for rotation is viable, and once selected that the entity can be adjusted to fit the semantics of that location. To do so, a line is drawn from the sensor to the entity’s center at the rotated position. The non-road points are voxelized and if the line intersects with any of these voxels, the mutation is not performed. If valid, the entity is aligned with the road at that position.

IV. EVALUATION

To evaluate our test generation technique for LiDAR perception systems, we pose the following research questions:

- **RQ1.** How effective is the technique at uncovering failures defined at different levels of severity and which mutations are the most effective?
- **RQ2.** How effective is the technique at generating realistic test cases?
- **RQ3.** How efficient is the technique in terms of the time to generate the mutations?

A. SUTs

We evaluate our technique through testing five highly competitive perception systems submitted to the SemanticKITTI benchmark for the “Semantic Segmentation” task [2], shown in Table II. We evaluate against these models as they were the top performing models that included code and a pre-trained model.

TABLE II: SUT Rankings

Rank		Name	Year	mIoU
Overall	With Code			
3	1	Cylinder3D [59]	2020	67.8
4	2	SPVNAS [43]	2020	66.4
5	3	JS3C-Net [52]	2021	59.5
11	5	SalsaNext [12]	2020	59.5
15	7	SqueezeSegV3 [51]	2021	55.9

Each model’s code was forked, edited to produce a consistent output format, and containerized in Docker for reproducibility.

B. Parameterization for Test Generation

We used the SemanticKITTI [2] data set as the backing data set for our mutations. As described in Section II-B, the initial Resource Collection stage identifies additional entities and then filters out entities not suitable for mutation. During this phase, the 259,347 entities in the data set and 8,443 additional entities identified were filtered to a set of 63,323 entities suitable for mutation³. We explore seven distinct compositions of the basic mutation operations, described in Table III, so selected to showcase the ability of the operations to produce a diverse set of test cases. These compositions and parameterizations are not exhaustive; e.g. the Add operation could be combined with the Scale, Deform, or Adjust Intensity operations to produce more variety in the added entities. Future work should examine directions for expansion that will provide meaningful utility in the types of tests generated.

TABLE III: Mutation Compositions Studied

Composition	Entity Classes Considered
Add ◦ Rotate	Any
Add ◦ Mirror ◦ Rotate	Any
Remove	Any
Adjust Intensity	Vehicle
Deform	Vehicle
Scale	Vehicle
Class-Preserving Swap	Sign

We generated 10,000 test cases for each of the 7 composite mutations. Full details of the values selected for the parameters described in Section III are available in the repository. The values selected represent reasonable defaults based on small-batch experimentation during tool design. For each test case, except for Remove mutations, the entity and mutation parameters, such as the location to deform or amount to adjust intensity, were chosen at random. Removal’s stringent invariants to ensure realism made it difficult for the technique to find suitable replacement points in *PC*, leading us to enumerate all 6,852 possible mutations instead of using random selection.

C. Results

1) *RQ1 Results: Failures Found at Different Thresholds:* As discussed in Section II-D, we investigate both the Accuracy and Jaccard metrics to assess an SUT’s performance. Since

the evaluation metric (Equation 2) proposed uses a tunable threshold, ϵ , to determine if an SUT passed or failed, we examine the effects of this parameter, grouping the failures observed by their drop in score; larger values of ϵ correspond to more severe failures. We consider ϵ values $\{(1, 2], (2, 3], (3, 4], (4, 5], (5, 100]\}$. Under the Accuracy metric, $\epsilon = 5$ corresponds to 5% of the PC being misclassified, indicating large portions of the scene are incorrect. Under the Jaccard metric, since there are 19 classes evaluated, this roughly corresponds to one class being entirely misclassified ($1/19 \approx 5.3\%$).

Fig. 5 shows the number of failures per SUT for each mutation type on the x-axis in log scale and the failure severity (as per the ϵ) on the y-axis; the left and right plots show the Accuracy and Jaccard metrics respectively. This shows that the tests generated by the technique identified between 17 and 101 failures among the 5 SUTs at the highest severity level ($\epsilon = 5$) under the Accuracy metric and at least 301 under the Jaccard metric. Breaking down the performance by SUT, we also see that JS3C-Net is considerably outperformed by the other SUTs, including those ranked lower in the original leaderboard shown in Table II.

To investigate further, we examine the breakdown of SUT performance based on the mutation type. Fig. 6 shows the number of failures at the highest level of severity ($\epsilon = 5$) per mutation. Every mutation yielded at least one failure except for the Sign Replace mutation under the Accuracy metric. As expected, the Jaccard metric identifies more failures, and more failures at higher thresholds, than the Accuracy metric since it is more sensitive to failures affecting small classes. Other than the Add Rotate and Add Mirror Rotate mutations, which yield similar failure results for all SUTs, the different SUTs are susceptible to different mutations. For example, JS3C-Net appears to be particularly brittle to Vehicle Deform and Vehicle Scale and to a lesser extent Remove and Vehicle Intensity.

We illustrate the type of failures found through Fig. 7 (others are provided in the repository), which shows a Vehicle Scale test case that uncovered a high-severity failure in JS3C-Net with an Accuracy drop of 21.9%. Fig. 7a shows JS3C-Net’s prediction on the unedited PC. The mutation then adjusts the scale of the vehicle on the right of the PC, and Fig. 7b shows the prediction on the mutated PC. Although the performance around the altered vehicle changes slightly, the accuracy drastically changes due to the large change in classification on the left side of the PC. In these images, orange corresponds to “fence” and yellow to “building”; as seen in Fig. 7c, the expected classification is fence for the entire area. While JS3C-Net classified more as fence than building before the mutation, after the mutation this reversed to more classified as building than fence, leading to a large decrease in accuracy. Such findings point to directions to improve robustness.

Fig. 8 illustrates another failure, this time in Cylinder3D identified by a drop in the Jaccard metric of 7.6%. This mutation added a bicyclist down the road from the LiDAR. Fig. 8a shows Cylinder3D’s actual prediction for this test case while Fig. 8b shows the mutated predicted semantics,

³A breakdown per class is available in our GitHub repository.

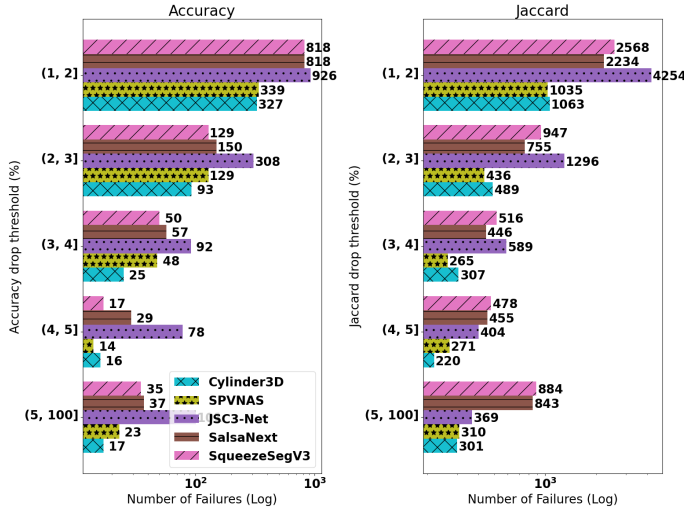


Fig. 5: Failure Counts at Different Thresholds

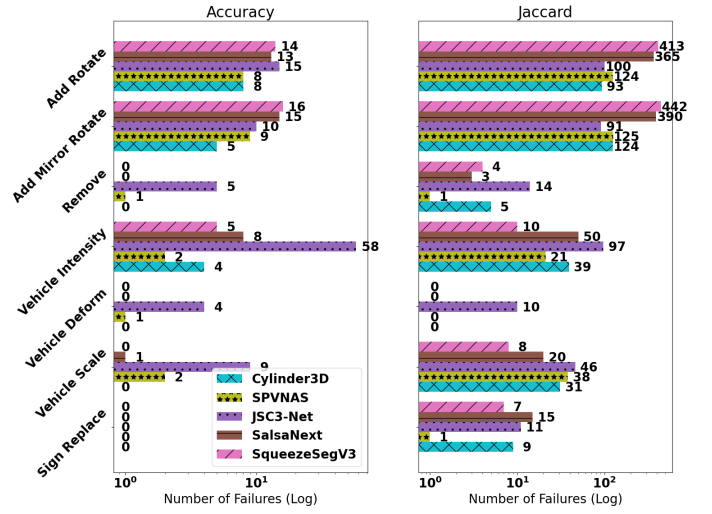


Fig. 6: SUT Failures by Mutation at Highest Severity

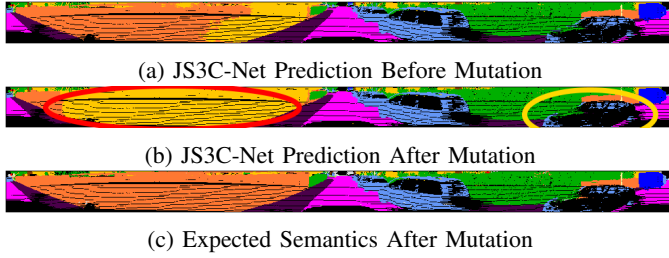


Fig. 7: JS3C-Net Performance on Vehicle Scale Mutation Prediction change in red oval, mutation change in yellow oval.

i.e. the expected output for Cylinder3D on the mutated input. Since in the SemanticKITTI coloring scheme the coloring for “bicyclist” is very similar to “road”, we have added a white outline in the picture to facilitate its visualization. Note that in Fig. 8a the bicyclist is composed of two regions, one cyan and one red corresponding to “bicycle” and “person” respectively. Thus, while Cylinder3D did identify the components of the “bicyclist” correctly, it did not identify them as a cohesive unit. Subtle differences such as this can be very important for perception and any components that rely on perception data. Under the SemanticKITTI labeling instructions, bicycles *cannot* be mobile while bicyclists can; this difference could be imperative to planning a safe trajectory for the system. While this difference is a 0.03% drop in the Accuracy metric, it is a 7.6% drop in the Jaccard metric due to the change in classes—“bicyclist” to “bicycle” and “person”—which Jaccard is uniquely suited to identify.

2) *RQ2 Results: Realism & False Positives:* When evaluating the performance of an SUT on a given test case, we want to ensure that the test case is realistic as a failing but unrealistic test case is a false positive. Recall from Section II-C that a test case is realistic if it could come from some possible configuration of the real world; thus, a failure on an unrealistic test case likely does not generalize to AS performance in the field—a

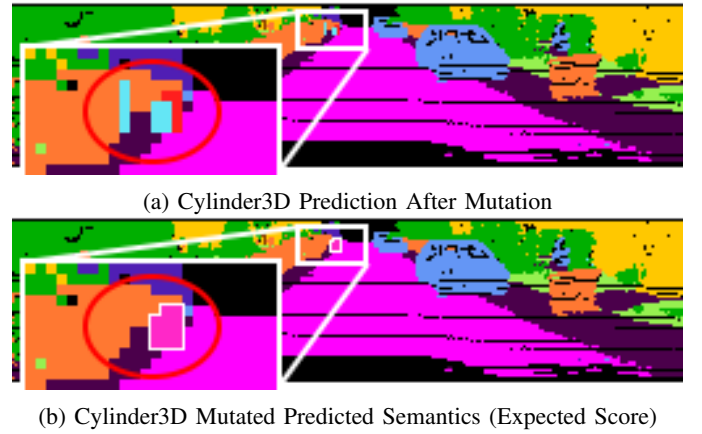


Fig. 8: Cylinder3D Performance on Add Rotate Mutation Prediction change in red oval.

TABLE IV: False Positive Counts, $V=3$

Mutation	Accuracy		Jaccard	
	TP	FP	TP	FP
Add Rotate	39	2	634	77
Add Mirror Rotate	42	2	672	79
Scene Remove	6	—	26	—
Vehicle Intensity	63	2	178	4
Vehicle Deform	5	—	10	—
Vehicle Scale	12	—	81	14
Sign Replace	—	—	37	1
Total (%)	167 (97%)	6 (3%)	1638 (90%)	175 (10%)

false positive⁴. The invariants described in Section II-C aim to prevent such false positives; however, these invariants must balance the need for realism with the ability to generate a diverse test set and in practice some false positives will occur.

⁴Unrealistic, but passing, tests are not considered. Our technique demonstrates utility by uncovering failures that could lead to unsafe behaviors; a passing test, realistic or not, only demonstrates safety on that specific input.

Determining if a given test case is realistic is a difficult and open problem. Realism exists on a continuum across multiple dimensions and measuring and codifying boundaries along this space remains an open challenge. Prior work on semantic mutations for images relied on human intuition to judge if an image was unrealistic, meaning that it could not have come from some configuration of the real world, and thus would be a false positive test failure [50]; however, this is subjective and cannot be automated. Further, human intuition is not suited to judging conformity of LiDAR PCs, as the sensor modality is far removed from how we perceive the world⁵. Related work on LiDAR input generation did not evaluate realism for individual test cases, instead using aggregate performance metrics to show similar results under real and generated test cases to argue that all test cases should be assumed to be realistic and thus not false positives [15], [30].

To address these limitations, we introduce a voting mechanism to determine if a generated input for a LiDAR perception system is realistic or not. Given a set of SUTs, S , trained on real world data, and a test case t' , we evaluate each $s \in S$ on t' . We then use the count of SUTs that fail on t' , v , to judge realism based on a voting threshold V . If $v < V$ SUTs fail, then we determine that the test case was realistic and thus not a false positive; i.e., the failure can be attributed to a fault in the SUT(s). If $v \geq V$ SUTs fail, then we determine that the failure is due to an unrealistic input and is a false positive; i.e., the SUTs failed because the input was not from a possible configuration of the real world, leading to a false positive. This method allows for *automated* false positive detection, but requires a diverse set of well-performing SUTs. Diversity is required as a realistic test case that results in failures for multiple similar systems will be judged as a false positive, leading to an overapproximation of false positives.

Table IV shows the number of false positive test cases for each of the mutation types. This includes only the highest level of failure examined in Section IV-C1 ($\epsilon = 5$). Selecting $V = 3$, corresponding to the majority of SUTs, we find only 3% and 10% of tests are deemed to be false positives under the Accuracy and Jaccard metrics respectively. Under the strictest criteria of $V = 2$, the rate increases to 19% and 36%, while under the most lenient criteria of $V = 5$ the rates drop to 0% and 1% respectively⁶. With the false positivity rate affected by both the choice of invariant parameters and V , this rate indicates a reasonable balance of finding impactful test cases while not spending too much execution time on false positives.

3) *RQ3 Results: Efficiency*: We now investigate the efficiency of the technique in terms of the time taken to generate the mutations. The average times per mutation were: 1.90s for Add Rotate, 1.97s for Add Mirror Rotate, 9.62s for Remove, 0.46s for Vehicle Intensity, 0.43s for Vehicle Deform, 0.84s for Vehicle Scale, and 0.71s for Sign Replace. The high time per test for Remove is due to the stringent invariant requirements. Overall, these mutations are generally slower

than the majority of the SUTs' time to label a scene—on average Cylinder3D took 0.59s, SPVNAS took 0.17s, JS3C-Net took 4.23s, SalsaNext took 0.03s, and SqueezeSegV3 took 0.16s. However, the technique test generation times are within one order of magnitude of test execution times. This is encouraging for several reasons. First, as a test generation technique, the tests created can be reused for future tests without incurring the full cost of generation, mitigating a high up-front cost. Second, although the costs of generation may be higher than the costs of evaluation, they are still pragmatic. Third, as the time taken to generate a mutation increases with the strictness of the invariants imposed, the timing is indirectly a tunable parameter and could be altered based on testing goals. Finally, as a research prototype, general efficiency refinements to the tool implementation have the opportunity to decrease the time per mutation.

D. Threats to Validity

The external validity of our findings are affected by our choice of data set and SUTs. We selected the SemanticKITTI data set for its popularity as a benchmark for LiDAR-based perception systems, and we selected the top five performing SUTs that were publicly available. Extending to additional available benchmarks and SUTs would help to generalize the findings. The internal validity of our findings may be affected by the implementation of the mutation pipeline which is complex and contains external components and many parameters. To mitigate these issues, we publicly release our code and provide documentation of the parameters used in the study, as well as all raw study results. Additionally, while the parameters chosen represent reasonable default values determined during system design, individual testing goals or testing using different data sets may necessitate changes to these values. In terms of construct validity, though we present a novel method for automatically assessing PC realism, the appropriate parameterization of this approach and other methods for external validation remain an ongoing challenge.

V. RELATED WORK

In this section, we explore the current state of testing LiDAR-based perception systems, focusing on the current weaknesses and differences from our approach.

A. Reality and Simulation

Real world testing is the gold standard for validation of ASs as it reflects the operating conditions the system will encounter when deployed; however, its effectiveness is constrained due to limited diversity. Although it is straight-forward to collect typical data in the real world by operating the system, it is impractical or dangerous to collect the safety-critical scenarios required for testing [29]. Further, labeling collected data with ground truth interpretation comes at great expense—the SemanticKITTI data set required over 1700 hours of human effort to annotate [2], and other data sets have noted the need for specialized training and software with multiple stages of

⁵Examples of false positive PCs as determined by our technique and further discussion on human detection are available in our GitHub repository.

⁶Full discussion of all choices of V is available in the repository.

review [42]. Combined, these factors limit the availability of diverse test inputs gathered from the real world.

To improve diversity, testing using simulation is widely used as new, unseen scenarios for testing can be generated relatively cheaply [34], [40], [45]. Existing approaches have used general purpose autonomous driving simulators [13], LiDAR-specific simulators [21], or adapted game engines [47], [56]. Prior work has noted that the major limitation of the simulation testing approach, the simulation-reality gap [27], remains a relevant issue in the LiDAR testing domain [47]. Recent work has attempted to improve realism by supplementing simulated data with real-world data [30], [48]; however, the simulation-reality gap remains a significant challenge.

B. Input Generation

In addition to real data and simulation-based methods, prior work has explored several lines of direct input generation: weather simulation, adversarial testing, and data augmentation. While each of these is an important part of LiDAR-based testing, we briefly discuss the limitations and differences with respect to our technique’s goal of diversity and realism.

1) *Weather Simulation*: LiDAR perception is affected by weather which can increase noise and occlude points [39]. Prior work has explored generating weather effects including snow, fog, and rain [22], [24], [28]. While validating weather performance is crucial to safety, the expected output in these weather-based test cases is the same as the expected output before weather effects are introduced, limiting diversity.

2) *Adversarial Testing*: Machine learned components are known to exhibit unstable behavior where a slight change in the input can lead to a drastic change in the output [20] and LiDAR-based perception systems are no exception [58]. This leads to adversarial attacks which try to alter the prediction in a specific way through small changes to the input, with the potential to produce system-level failures [7], [46], [60]. While improving robustness is important to AS safety, the niche inputs do not improve general test suite diversity.

3) *Data Augmentation*: Data augmentation techniques take as input a prior input and expected output pair and generate a novel input and expected output pair. Several techniques for images have been extended to LiDAR PCs; for example, affine transformations [22], or fusing together different scenes [26]. Due to this global effect, the semantics of each point is the same before and after even if the point is moved, limiting improvements to diversity. Prior entity-specific augmentations have removed points, added noise, or mixed together two entities [8], though without considerations for realism. Another technique for augmentation randomly took points from other scenes and added them to the input scene [41], [53], [55]; however, the randomness could lead to physically infeasible point locations. Following techniques improved this by using 3D models and simulating how the object’s addition would affect the point locations [15]. Although this improves diversity, the lack of consideration for PC intensity may affect realism. Recent literature explored the insertion of objects from one scene into another scene to improve training diversity [26].

While this improves diversity, the entity mutations considered are limited to whole-entity changes such as mirroring and do not consider localized mutations such as deformations. Further, as a technique focused on training diversification, it lacks a mechanism for generating test oracles.

VI. CONCLUSION AND BROADER IMPACT

Our approach provides a first-of-its-kind semantic mutation testing technique for failure detection in LiDAR-based AS perception systems using realistic data to explore never-before-seen scenarios. It provides a transition from simulation testing to real-world testing. Whereas initial research and development for ASs likely occurs in simulation, allowing rapid iteration, the simulation-reality gap necessitates testing on real-world data which involves ethical considerations, steep costs to execute tests, and steeper costs for failures. Our technique bridges this gap, allowing for increased testing on difficult-to-obtain real-world data before the system is exercised or deployed in the real world with the potential to save time, cost, and lives by identifying safety-critical failures earlier. Our implementation and its usage as part of a study demonstrates the ability of the approach to efficiently find potentially safety-critical failures in LiDAR-based perception systems. The approach can be naturally extended to accommodate mutations for other domains such as urban UAVs, indoor robotics, etc.—the same add, remove, change paradigm and realism invariant framework apply to all domains, though specific domains may have additional change operations and accompanying realism invariants, e.g. when targeting UAV perception the realism invariants related to added entities being on the ground may be replaced with invariants about an entity’s possible altitude.

This sets several directions for future work. Although we demonstrate the ability for random selection to find failures, future work should explore feedback-driven, guided parameter selection. Prior literature has shown the effectiveness of using SUT performance to guide testing in a robotic context [17], [33], [49]. Additionally, ASs use LiDAR in conjunction with other sensors simultaneously and may integrate multiple sensors for perception; future work should explore mixed-input generation that mutates data across multiple sensor modalities simultaneously. Further, ASs operate in continuous time and may use prior perception knowledge when perceiving the current environment; future work should explore mutations over sequences of point clouds and explore the richer mutation space available across multiple time steps.

VII. DATA AVAILABILITY

The data, artifacts, and additional discussion, much of which we could not include in the paper, are available in the repository on GitHub and archived on Software Heritage [10].

ACKNOWLEDGEMENTS

This work was supported in part by funds provided by NSF#1924777, NSF#1909414, and AFOSR#FA9550-21-1-0164. Trey Woodlief was supported by a University of Virginia SEAS Fellowship.

REFERENCES

- [1] Evan Ackerman. Ford and baidu invest \$150 million in velodyne for affordable lidar for self-driving cars. *IEEE Spectrum*, Aug 2016.
- [2] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences. In *Proc. of the IEEE/CVF International Conf. on Computer Vision (ICCV)*, 2019.
- [3] Rebecca Bellan. Motional launches autonomous hyundai ioniq 5s on lyft network in las vegas. *TechCrunch*, Aug 2022.
- [4] NTS Board. Collision between vehicle controlled by developmental automated driving system and pedestrian. nat. transpot. saf. board, washington, dc. Technical report, USA, Tech. Rep. HAR-19-03, 2019. URL <https://www.nts.gov/investigations...>, 2019.
- [5] Neal E Boudette and Niraj Chokshi. U.s. will investigate tesla's autopilot system over crashes with emergency vehicles. *New York Times*, Aug 2021.
- [6] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020.
- [7] Yulong Cao, Chaowei Xiao, Dawei Yang, Jing Fang, Ruigang Yang, Mingyan Liu, and Bo Li. Adversarial objects against lidar-based autonomous driving systems. *CoRR*, abs/1907.05418, 2019.
- [8] Jaeseok Choi, Yeji Song, and Nojun Kwak. Part-aware data augmentation for 3d object detection in point cloud. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3391–3397, 2021.
- [9] Yookyung Choi, Namil Kim, Soonmin Hwang, Kibaek Park, Jae Shin Yoon, Kyoungwan An, and In So Kweon. Kaist multi-spectral day/night data set for autonomous and assisted driving. *IEEE Transactions on Intelligent Transportation Systems*, 19(3):934–948, 2018.
- [10] Garrett Christian, Trey Woodlief, and Sebastian Elbaum. semLidarFuzz. <https://archive.softwareheritage.org/whl:1:dir:a77c06955132e4b2d763184476fba2694c9c331b;origin=https://github.com/less-lab-uva/semLidarFuzz;visit=swl:1:snp:48b0bfc51ff2317f240bb44037fc20ecb0a3f79;anchor=swl:1:rev:c389f46acb46a17a5c2b8a75feaf7ebe19d41071>, 2023.
- [11] Daniel Cohen-Or and Arie Kaufman. Fundamentals of surface voxelization. *Graphical models and image processing*, 57(6):453–461, 1995.
- [12] Tiago Cortinhal, George Tzelepis, and Eren Erdal Aksoy. Salsanext: Fast semantic segmentation of lidar point clouds for autonomous driving. *CoRR*, abs/2003.03653, 2020.
- [13] Alexey Dosovitskiy, Germán Ros, Felipe Codevilla, Antonio M. López, and Vladlen Koltun. CARLA: an open urban driving simulator. *CoRR*, abs/1711.03938, 2017.
- [14] Elmar Eisemann and Xavier Décoret. Fast scene voxelization and applications. In *Proceedings of the 2006 symposium on Interactive 3D graphics and games*, pages 71–78, 2006.
- [15] Jin Fang, Xinxin Zuo, Dingfu Zhou, Shengze Jin, Sen Wang, and Liangjun Zhang. Lidar-aug: A general rendering-based augmentation framework for 3d object detection. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4708–4718, 2021.
- [16] Ryan Felton. Lidar maker velodyne shifts away blame in fatal uber self-driving crash, Mar 2018.
- [17] Alessio Gambi, Marc Mueller, and Gordon Fraser. Automatically testing self-driving cars with search-based procedural content generation. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 318–328, 2019.
- [18] Joshua Garcia, Yang Feng, Junjie Shen, Sumaya Almanee, Yuan Xia, and Qi Alfred Chen. A comprehensive study of autonomous vehicle bugs. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*, pages 385–396, 2020.
- [19] Richard Gonzales. Feds say self-driving uber suv did not recognize jaywalking pedestrian in fatal crash, Nov 2019.
- [20] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [21] Michael Gschwandtner, Roland Kwitt, Andreas Uhl, and Wolfgang Pree. Bensor: Blender sensor simulation toolbox. In George Bebis, Richard Boyle, Bahram Parvin, Darko Koracin, Song Wang, Kim Kyungnam, Bedrich Benes, Kenneth Moreland, Christoph Borst, Stephen DiVerdi, Chiang Yi-Jen, and Jiang Ming, editors, *Advances in Visual Computing*, pages 199–208, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [22] An Guo, Yang Feng, and Zhenyu Chen. Lirtest: Augmenting lidar point clouds for automated testing of autonomous driving systems. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2022*, page 480–492, New York, NY, USA, 2022. Association for Computing Machinery.
- [23] T Hackel, N Savinov, L Ladicky, JD Wegner, K Schindler, and M Pollefeys. Semantic3d. net: a new large-scale point cloud classification benchmark. *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences*, 41:91–98, 2017.
- [24] Martin Hahner, Christos Sakaridis, Mario Bijelic, Felix Heide, Fisher Yu, Dengxin Dai, and Luc Van Gool. Lidar snowfall simulation for robust 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16364–16374, June 2022.
- [25] Charles R. Harris, K. Jarrod Millman, Stéfán J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020.
- [26] Frederik Hasecke, Martin Alsasser, and Anton Kummert. What can be seen is what you get: Structure aware point cloud augmentation. 2022.
- [27] Nick Jakobi, Phil Husbands, and Inman Harvey. Noise and the reality gap: The use of simulation in evolutionary robotics. In *European Conference on Artificial Life*, pages 704–720. Springer, 1995.
- [28] Velat Kilic, Deepti Hegde, Vishwanath Sindagi, A. Brinton Cooper, Mark A. Foster, and Vishal M. Patel. Lidar light scattering augmentation (lisa): Physics-based simulation of adverse weather conditions for 3d object detection, 2021.
- [29] Waymo LLC. Waymo's safety methodologies and safety readiness determinations. Technical report, October 2020.
- [30] Sivabalan Manivasagam, Shenlong Wang, Kelvin Wong, Wenyuan Zeng, Mikita Sazanovich, Shuhan Tan, Bin Yang, Wei-Chiu Ma, and Raquel Urtasun. Lidsim: Realistic lidar simulation by leveraging the real world, 2020.
- [31] Post author By Luminar Marketing. Luminar and volvo cars announce plans for next generation suv to be revealed in 2022, Jan 2022.
- [32] Aarian Marshall. Uber video shows the kind of crash self-driving cars are made to avoid, Mar 2018.
- [33] Cu D Nguyen, Simon Miles, Anna Perini, Paolo Tonella, Mark Harman, and Michael Luck. Evolutionary testing of autonomous software agents. *Autonomous Agents and Multi-Agent Systems*, 25:260–283, 2012.
- [34] Matthew O'Kelly, Aman Sinha, Hongseok Namkoong, Russ Tedrake, and John C Duchi. Scalable end-to-end autonomous vehicle testing via rare-event simulation. *Advances in neural information processing systems*, 31, 2018.
- [35] Abhishek Patil, Srikanth Malla, Haiming Gang, and Yi-Ting Chen. The h3d dataset for full-surround 3d multi-object detection and tracking in crowded urban scenes. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 9552–9557. IEEE, 2019.
- [36] Jake Port. From the vault: Lidar: How self-driving cars 'see' where they're going, Jan 2022.
- [37] Sabbir Rangwala. Amazon's zoox acquisition– is lidar next? *Forbes*, Aug 2020.
- [38] Xavier Roynard, Jean-Emmanuel Deschaud, and François Goulette. Paris-lille-3d: A large and high-quality ground-truth urban point cloud dataset for automatic segmentation and classification. *The International Journal of Robotics Research*, 37(6):545–557, 2018.
- [39] Fredrik Schalling, Sebastian Ljungberg, and Naveen Mohan. Benchmarking lidar sensors for development and evaluation of automotive perception. In *2019 4th International Conference and Workshops on Recent Advances and Innovations in Engineering (ICRAIE)*, pages 1–6, 2019.
- [40] Hans-Peter Schöner. Simulation in development and testing of autonomous vehicles. In *18. Internationale Stuttgarter Symposium*, pages 1083–1095. Springer, 2018.
- [41] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointcnn: 3d object proposal generation and detection from point cloud. *CoRR*, abs/1812.04244, 2018.

- [42] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020.
- [43] Haotian Tang, Zhijian Liu, Shengyu Zhao, Yujun Lin, Ji Lin, Hanrui Wang, and Song Han. Searching efficient 3d architectures with sparse point-voxel convolution. *CoRR*, abs/2007.16100, 2020.
- [44] Brad Templeton. Tesla in taiwan crashes directly into overturned truck, ignores pedestrian, with autopilot on. *Forbes*, Jun 2020.
- [45] Christopher Steven Timperley, Afsoon Afzal, Deborah S Katz, Jam Marcos Hernandez, and Claire Le Goues. Crashing simulated planes is cheap: Can simulation detect robotics bugs early? In *2018 IEEE 11th International Conference on Software Testing, Verification and Validation (ICST)*, pages 331–342. IEEE, 2018.
- [46] James Tu, Mengye Ren, Siva Manivasagam, Ming Liang, Bin Yang, Richard Du, Frank Cheng, and Raquel Urtasun. Physically realizable adversarial examples for lidar object detection. *CoRR*, abs/2004.00543, 2020.
- [47] Patrik Vacek, Otakar Jašek, Karel Zimmermann, and Tomáš Svoboda. Learning to predict lidar intensities. *IEEE Transactions on Intelligent Transportation Systems*, 23(4):3556–3564, 2022.
- [48] Jingkan Wang, Ava Pun, James Tu, Sivabalan Manivasagam, Abbas Sadat, Sergio Casas, Mengye Ren, and Raquel Urtasun. Advsim: Generating safety-critical scenarios for self-driving vehicles, 2021.
- [49] Trey Woodlief, Sebastian Elbaum, and Kevin Sullivan. Fuzzing mobile robot environments for fast automated crash detection. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5417–5423. IEEE, 2021.
- [50] Trey Woodlief, Sebastian Elbaum, and Kevin Sullivan. Semantic image fuzzing of ai perception systems. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*, pages 1958–1969, 2022.
- [51] Chenfeng Xu, Bichen Wu, Zining Wang, Wei Zhan, Péter Vajda, Kurt Keutzer, and Masayoshi Tomizuka. Squeezesegv3: Spatially-adaptive convolution for efficient point-cloud segmentation. In *ECCV*, 2020.
- [52] Xu Yan, Jiantao Gao, Jie Li, Ruimao Zhang, Zhen Li, Rui Huang, and Shuguang Cui. Sparse single sweep lidar point cloud segmentation via learning contextual shape priors from scene completion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3101–3109, 2021.
- [53] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10), 2018.
- [54] Yingzhi Yang, Zhang Yan, and Brenda Goh. Baidu ceo says ev arm’s autonomous driving tech will be ahead of tesla. *Reuters*, Aug 2022.
- [55] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Center-based 3d object detection and tracking. *CoRR*, abs/2006.11275, 2020.
- [56] Xiangyu Yue, Bichen Wu, Sanjit A. Seshia, Kurt Keutzer, and Alberto L. Sangiovanni-Vincentelli. A lidar point cloud generator: from a virtual world to autonomous driving, 2018.
- [57] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018.
- [58] Zhi Quan Zhou and Liqun Sun. Metamorphic testing of driverless cars. *Commun. ACM*, 62(3):61–67, feb 2019.
- [59] Xinge Zhu, Hui Zhou, Tai Wang, Fangzhou Hong, Yuexin Ma, Wei Li, Hongsheng Li, and Dahua Lin. Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. *CoRR*, abs/2011.10033, 2020.
- [60] Yi Zhu, Chenglin Miao, Foad Hajiaghajani, Mengdi Huai, Lu Su, and Chunming Qiao. Adversarial attacks against lidar semantic segmentation in autonomous driving. In *Proceedings of the 19th ACM Conference on Embedded Networked Sensor Systems, SenSys ’21*, page 329–342, New York, NY, USA, 2021. Association for Computing Machinery.