# Property-Based Testing in Simulation for Verifying Robot Action Execution in Tabletop Manipulation

Salman Omar Sohail[†], Alex Mitrevski[†§], Nico Hochgeschwender[†], and Paul G. Plöger[†]

*Abstract*— An important prerequisite for the reliability and robustness of a service robot is ensuring the robot's correct behavior when it performs various tasks of interest. Extensive testing is one established approach for ensuring behavioural correctness; this becomes even more important with the integration of learning-based methods into robot software architectures, as there are often no theoretical guarantees about the performance of such methods in varying scenarios. In this paper, we aim towards evaluating the correctness of robot behaviors in tabletop manipulation through automatic generation of simulated test scenarios in which a robot assesses its performance using property-based testing. In particular, key properties of interest for various robot actions are encoded in an action ontology and are then verified and validated within a simulated environment. We evaluate our framework with a Toyota Human Support Robot (HSR) which is tested in a Gazebo simulation. We show that our framework can correctly and consistently identify various failed actions in a variety of randomised tabletop manipulation scenarios, in addition to providing deeper insights into the type and location of failures for each designed property.

## I. INTRODUCTION

With the integration of autonomous service robots into industries and households, there is a requirement for increased robot safety and dependability. To fulfill these requirements, robot developers need to validate their systems extensively before deployment. In principle, validation is a challenging problem due to factors such as the environment's unpredictability, the robot's lack of knowledge, hardware and software failures, but also due to the fact that a robot may utilise learning-based components for which formal correctness guarantees are difficult to provide.

One common approach for robot verification and validation is testing directly in the real world, but real-world tests are often complex to set up and perform to an extent that would provide sufficient test coverage of a complete robot system. An alternative to this is simulation-based testing, which provides multiple attractive features, such as comparatively low setup and execution costs, speed, scalability, robot safety, as well as a possibility to automate tests [1], at the cost of sacrificing the realism of real-world testing. Setting up and executing simulated tests can, however, be a challenging problem on its own. In particular, simulated scenarios for a robot are often hand-crafted and each scenario
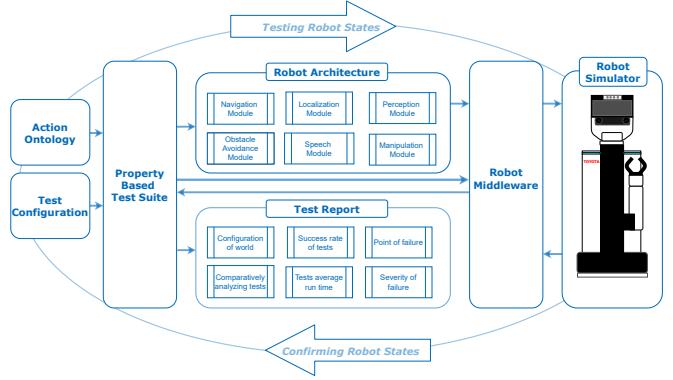


Fig. 1: Our proposed simulation-based property-based testing framework for verifying robot action execution

has to be manually altered when the need arises to test different scenarios [2].[1] In addition, exhaustive robot tests need to cover not only the software aspects of the robot system, but also the overall behavior and decisions that the robot makes during task execution [4].

This paper is motivated by these challenges in robot testing and addresses the question of how to test robot behaviours so that robot execution failures, which go beyond software failures, can be systematically analysed. Our proposed framework uses automatic scenario generation for simulation-based robot testing, namely we generate a set of simulated scenarios in which a robot performs various actions and assesses its performance using property-based testing [5]. The framework includes four core components: (i) an action ontology that relates primitive robot actions to properties of the world that determine the success or failure when those actions are performed, (ii) scenario generation based on which randomised world configurations are created to mimic situations that a robot might encounter, (iii) property-based tests that evaluate the performance of a robot in the generated scenarios, and (iv) report generation that includes the world properties used in the scenario generation as well as the results of the executed property tests, both of which can be used to identify potential reasons for an execution failure. An overview of our framework is provided in Fig. 1. We demonstrate the use of the framework using a Toyota Human Support Robot (HSR) [6] and our software stack for domestic robots, focusing on functionalities such as point-based navigation, object perception for subsequent manipulation, as well as object pickup and placing in various

[†]The authors are with the Autonomous Systems Group, Department of Computer Science, Hochschule Bonn-Rhein-Sieg, Sankt Augustin, Germany `salman.sohail@smail.inf.h-brs.de`, `<aleksandar.mitrevski, nico.hochgeschwender, paul.ploeger>@h-brs.de`

[§]Corresponding author

---

[1]For instance, this has been the case in some of our earlier work [3].

common domestic object manipulation scenarios.

## II. RELATED WORK

In the context of pure software systems, various testing strategies can be used, such as model-based, functional, mutation, or regression testing [7]; however, as pointed out by Bihlmaier and Wörn [8], software testing alone cannot capture the full complexity of robot systems, which interact with the world and other agents. Robot systems thus require dedicated testing methodologies that are able to incorporate physical aspects of the world and take into account the potential of experiencing execution failures.

In the context of simulation-based testing, the core aspect is the scenario generation process. Scenario generation has been studied fairly thoroughly for autonomous driving. Majumdar et al. [9] present a test specification language called PARACOSM, which creates realistic environment simulations for testing autonomous driving systems in the Unity physics engine; the language allows defining object and environment properties in a test scenario, thereby providing a declarative way of scenario specification and configuration. Park et al. [10] propose an approach for simulated scenario generation based on networks trained with real-world driving data; the objective here is to generate realistic driving scenarios that represent multiple vehicle events in a single video session, which enables learning data generation and subsequent vehicle testing. In [11], Bayesian Optimisation is used to generate stimulating adversarial scenarios for an autonomous vehicle; imitation learning is subsequently used to acquire corrective actions from an expert while the scenario is being executed in simulation. A similar technique is proposed by Koren et al. [12], where Adaptive Stress Testing based on a Markov Decision Process is used to identify driving scenarios that are likely to lead to a failure. Our framework draws upon various insights from the above work. As in [9], we use a specification language to control various test generation parameters. Similar to [10], our scenarios may represent multiple events, but we do not rely on available video data for training a test generator, although it would be possible to integrate the use of video data through an appropriate test generation strategy. Finally, methods similar to [11], [12] could be combined with property-based testing by learning to focus on challenging test scenarios, which provides an interesting direction for future work.

In the context of service robots, various testing methodologies have been used. In [13], property-based testing is used for testing robots whose software is based on the Robot Operating System (ROS), such that tests are based on properties about the expected input-output relations represented via the ROS component graph. Araiza et al. [4] compare two test generation strategies - Belief-Desire-Intention and model checking of timed automata - which can be used for testing collaborative service robots. In [14], Estivill-Castro et al. propose a framework for simulation-based robot testing that utilises continuous integration; here, tests are centered around robot behaviours which are modelled by automata, such that the test complexity is governed by a hierarchical structure underlying the behaviours. As [13], our work is based on property-based testing; however, while our robot software is also embedded in ROS, our general approach is independent of a communication framework, namely it allows testing behaviours that do not require a communication framework to be used during execution. Similar to [4], our test generation strategy is based on modelled relations between robot actions and properties of interest, but our primary focus is on physical interaction with the world and failures that may occur during action execution with a large variety of objects. As in [14], our framework allows robot testing at different levels of complexity, namely we aim to test both individual robot actions as well as complete tasks.

Robot simulations can also be utilised to guide the robot execution process at runtime, which can be seen as an online testing strategy. Kunze et al. [15] propose a simulation-based methodology to identify suitable parameters for executing a robot action by recording logs from simulated executions and analysing those using the event calculus. Similarly, Mösenlechner and Beetz [16] present a methodology to identify a set of valid positions and orientations for a mobile robot's base by considering predictions from simulated events while picking and placing objects. Our simulated tests are similar to the scenario executions in [15], [16], but our framework is only intended to be used for comprehensive offline testing of a robot's behavioural components so that failures can be analysed and improved before a robot is practically deployed.

## III. PROPERTY-BASED TESTING IN SIMULATION

Our proposed framework aims at testing robot actions, such as *pick object* or *move to*, and complete tasks that are composed of multiple actions, using property-based testing in simulation, such that the objective is to automate the test generation and evaluation process so that actions can be tested under different conditions and with different parameters. To make this possible, it is necessary to specify how test scenarios should be generated and which properties should be tested for each action. In this section, we start with a description of property-based testing and explain how we encode and test properties, as well as how tests are documented for subsequent analysis.

### A. Property-Based Robot Testing

Property-based testing is a software testing framework that verifies whether certain specifications, or properties, are met in a given software system [5], such that the idea is to feed a component with varied inputs and check whether it fails with any given input values.[2] Property-based testing is a more general version of unit testing, as a single property-based test can substitute a collection of manually-specified unit tests; in other words, property-based testing can be seen as generative testing that can cover a large amount of the domain space using parameter generation *strategies*.

---

[2]The mechanism of varied inputs fed into property-based tests is somewhat similar to fuzz testing in which random and possibly invalid data is passed into a software component to detect failures.
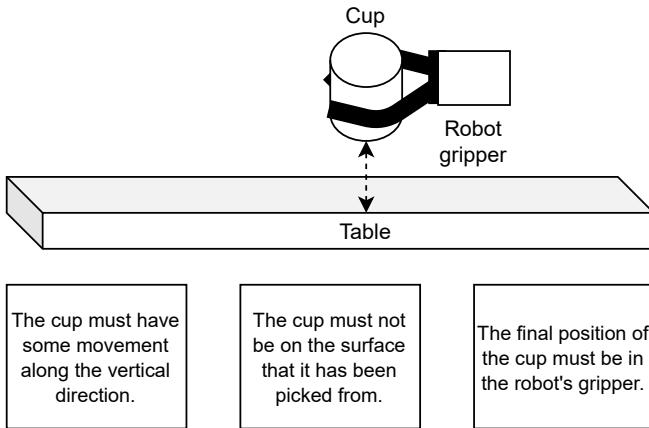
**Expected properties of a cup that has been picked-up**



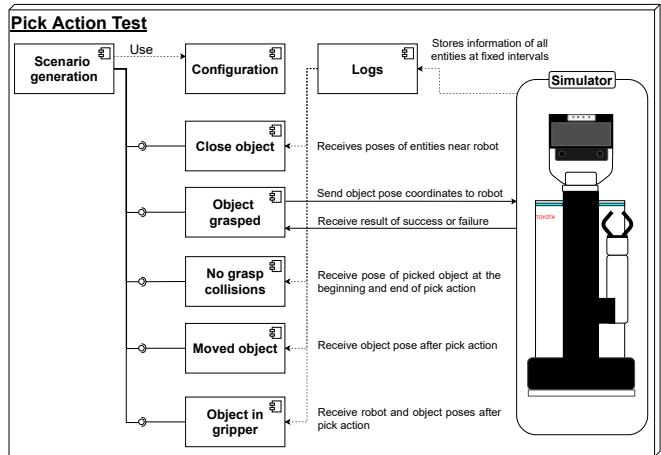Fig. 2: Relevant properties for an object grasping action



Fig. 3: An implementation example of modelling properties for a *pick object* action of a robot.

TABLE I: Defined properties for an object grasping action

| Property | Brief Description |
|---|---|
| Close object | Objects should exist within the proximity of the robot before the action |
| Object grasped | The action should report that it has completed its execution |
| No grasp collisions | The positions of all objects, with the exception of the robot and the object to be grasped, should remain the same before and after the action |
| Moved object | The position of the grasped object should change after the action is executed |
| Object in gripper | The grasped object should be in the robot's hand at the end of the action |

The general structure for applying property-based testing consists of (i) modelling the expected behavior of a system in terms of properties, (ii) determining the range of parameters that is of interest, (iii) activating an input generator in order to generate parameters in the specified range, and (iv) verifying that the test result satisfies the specified properties.

An example of applying property-based testing is given by the simple scenario shown in Fig. 2, in which we want to check whether a robot has picked up an object. The following properties are defined for an object to be picked up: (i) the object must be elevated from the surface that it was picked from, (ii) the object must not be on the surface that it was picked from, and (iii) the object's final position must be in the robot's hand. According to the first property, we expect that a picked up object will always be elevated from its original position on the surface; a violation of this property should result in a failure, as this would indicate that the object was knocked over while the robot was attempting to grasp it or the object was never grasped in the first place. Similarly, the second property identifies whether a successfully grasped object has slipped from the gripper back onto the table during or after the pick action. The third property verifies that the object is in the robot's hand. Through a combination of these three property tests, different types of failures can be identified; for instance, the success of property one and two and failure of property three would mean that the robot successfully picked up the object, but the object did not remain in its hand and fell down on the floor. Defining more properties for a system increases the testing mechanism's ability to identify the point of failure. To increase the coverage for a given use cases, we utilise *input generators*, which pass varied parameters for verifying the system's properties. As described later, these generators are parameterised at test runtime using a minimal configuration.

A more complete definition of relevant properties that we use for testing an object grasping action is given in Table I; the testing process itself is illustrated in Fig. 3.

Similar properties are also encoded for other actions involved in tabletop manipulation, in particular *move to*, *perceive plane*, and *place object*.

*B. Action Ontology*

We encode the information about which properties to test for a given action in an ontology $\mathcal{O}$ represented in the Web Ontology Language (OWL).[3] In the ontology, each action is represented as an instance of an `Action` class. Actions are performed with respect to a predefined coordinate frame and are associated with multiple properties that need to be tested for verifying the execution success. Properties are modelled as instances of a `Property` class, where each property depends on parameters that are set in the context of a given action; each such parameter is represented as an instance of a `PropertyParameter` class, such that the values of the parameters are set at runtime during test generation.

In the ontology, these classes are related to each other through the object properties shown in Table II. Here, the `performedIn` property relates an action to the coordinate frame in which the action is performed. `successProperty` specifies which properties need to be verified for the execution of an action to be considered successful. `needsParameter` and `hasParameter` specify which parameters are required by a given property

---

[3]https://www.w3.org/OWL/

TABLE II: Ontology properties

| Object property | Domain | Range |
|---|---|---|
| `performedIn` | `Action` | `CoordinateFrame` |
| `successProperty` | `Action` | `Property` |
| `hasParameter` | `Action` | `PropertyParameter` |
| `needsParameter` | `Property` | `PropertyParameter` |

**Algorithm 1** Property-based action testing

1: **function** TESTACTION($a$, $\mathcal{O}$)
2: $\quad P_a \leftarrow$ getActionParameters($\mathcal{O}$, $a$)
3: $\quad P_{val} \leftarrow \{\}$
4: $\quad$ **for** $p$ **in** $P_a$ **do**
5: $\quad\quad p_{val} \leftarrow$ generateParameter($p$)
6: $\quad\quad P_{val} \leftarrow P_{val} \cup p_{val}$
7: $\quad$ executeAction($P_{val}$)
8: $\quad Props_a \leftarrow$ getSuccessProperties($\mathcal{O}$, $a$)
9: $\quad$ **for** pp **in** $Props_a$ **do**
10: $\quad\quad pp_{val} \leftarrow$ getPropertyParameters(pp, $P_{val}$)
11: $\quad\quad$ **if not** pp($pp_{val}$) **then**
12: $\quad\quad\quad$ **return false**
13: $\quad$ **return true**

and which parameters are set for testing a given action, respectively. To indicate that the values of property parameters are assigned at runtime, we represent each instance of `PropertyParameter` by a designator, similar to [17]. An example ontology model of a *pick object* action, which illustrates what is encoded for each action, is shown below.[4]

Listing 1: OWL-based property model of an object grasping action

```
<owl:NamedIndividual rdf:about="http://actions#Pick">
    <rdf:type rdf:resource="http://actions#Action"/>

    <action:performedIn rdf:resource="http://actions#BaseLink"/>

    <action:successProperty rdf:resource="http://actions#MovedAlongY"/>
    <action:successProperty rdf:resource="http://actions#OnSurface"/>
    <action:successProperty rdf:resource="http://actions#InHand"/>

    <action:hasParameter rdf:resource="http://actions#objectLocation"/>
    <action:hasParameter rdf:resource="http://actions#objectSurface"/>
    <action:hasParameter rdf:resource="http://actions#pickObject"/>
    <action:hasParameter rdf:resource="http://actions#surfaceObjects"/>
</owl:NamedIndividual>

<owl:NamedIndividual rdf:about="http://actions#OnSurface">
    <rdf:type rdf:resource="http://actions#Property"/>
    <action:needsParameter rdf:resource="http://actions#objectSurface"/>
    <action:needsParameter rdf:resource="http://actions#pickObject"/>
</owl:NamedIndividual>
```

Automatic testing of properties is possible by implementing each property $pp$ as a parameterisable function, such that the properties $Props_a$ that need to be tested for a given action $a$ are dynamically invoked. Alg. 1 provides a high-level summary of the property testing procedure.

*C. Test Scenario Generation*

Similar to [18], we define a test scenario as a combination of an environment configuration and an objective that a robot needs to achieve. The scenario generation process thus requires to (i) set up an environment and (ii) assign a task to the robot being tested. Our scenario generation component serves as a parameterisable input generator which creates scenarios by placing objects in different positions and

orientations in a given environment. The generator uses a similar approach to [19]; in particular, to generate diverse and dynamic scenarios, we use custom 3D models of different household objects, such that scenarios are generated by placing the models in randomised poses.[5] In this paper, we focus on tabletop manipulation scenarios, so objects are generated on surfaces where the robot can manipulate them.

To allow for controlled test coverage, the scenario generation is configured through the parameters in Table III.

TABLE III: Scenario configuration parameters

| Parameter | Description |
|---|---|
| `tests` | List of tests to execute; each test is specified as a list of actions |
| `test_count` | Number of times to run each test |
| `test_launcher` | Path to a launch file that starts all components required by the tests |
| `model_dir` | Path to a directory of 3D models used in the tests |
| `worlds` | List of possible environments in which a test scenario can take place |
| `model_list` | Object models that can be used for manipulation |
| `nav_obstacle_list` | Object models that can be placed as navigation obstacles |
| `nav_obstacle_count` | Number of navigation obstacles to place in the environment |
| `location_list` | Names of locations to which the robot can navigate |
| `object_surfaces` | List of surfaces on which objects can be placed for manipulation |
| `place_object_surfaces` | List of possible surfaces to which objects can be brought |

These parameters are passed to the scenario generator through a TOML[6] configuration file. It should be noted that some of the configuration parameters are used to assign the values of property parameters specified in the ontology. This connection is established through TOML tables, where the header specifies the name of the property parameter that should be set and the corresponding configuration parameter defines the list of possible values. The format of the scenario configuration file for a *pick* action test is shown below.

Listing 2: Scenario generation configuration

```
tests = [["Pick"]]
test_count = int
test_launcher = str
model_dir = str
nav_obstacle_list = List[str]
nav_obstacle_count = int

[world]
worlds = List[str]

[objectLocation]
location_list = List[str]

[objectSurface]
object_surfaces = List[str]

[pickObject]
model_list = List[str]

[surfaceObjects]
model_list = List[str]
```

---

[5]To avoid models from spawning on top of each other due to the randomised nature of the placement, we apply 3D collision prevention using axis-aligned bounding boxes (AABB).

[6]https://github.com/toml-lang/toml

[4]The ontology is available at https://github.com/b-it-bots/action-execution/

## Failed  test_destination_verification

Overview    History    Retries

failed 14/12/2020 at 20:47:52

AssertionError: assert 0 <= (-0.8109999999999999 + 0.45)

failed 14/12/2020 at 20:47:07

AssertionError: assert [0, -2] <= (-0.8109999999999999 + 0.45)

broken 14/12/2020 at 20:46:04

IndexError: list index out of range

failed 14/12/2020 at 20:44:22

```
AssertionError: assert 5 <= (-0.8109999999999999 + 0.45)
 + where 5 = <tests.nav_test.TestNavigation instance at
0x7f9065159cf8>.coord_x
```

Fig. 4: A generated test report for four failed point-to-point navigation runs. In three of these, the navigation goal could not be reached due to a blocking obstacle, while one run failed due to a software error. *Overview* shows information about the overall performance of a test over multiple runs, *History* information about the number of runs of a test, and *Retries* details about the retries of a failed or broken test.

### D. Test Report Generation

For each test, we generate an automatic test report, which is an HTML file created using the *allure* library[7] in conjunction with *Hypothesis* [20] for property-based testing. The report is generated from json files that are created after each test run, such that it stores information about the test description and duration, world properties and parameters used for the tests, the test consistency and history, as well as the success rate and points of failure of the tests. Test behavior is described by the success rate and consistency of a test over multiple runs; in particular, for each test, the time taken is recorded, such that if the test is inconsistent over multiple runs, the test is considered unstable. Fig. 4 illustrates a test report for failed navigation tests.

### E. Test Suite Evaluation

We evaluate a test suite, which consists of multiple test runs, through the number of satisfied properties for each involved action over all runs. Let $T$ be a test suite in which $l$ different actions need to be tested. Actions are tested in $m$ randomly generated scenarios $S_k, 1 \le k \le m$ with different test parameters; each action $A_i^k, 1 \le i \le l$ has $n$ properties $pp_j^i, 1 \le j \le n$ that are evaluated by an indicator function $\mathbb{I}$. The evaluation of a test suite is thus measured through the normalised success of the individual scenarios:

$$T = \frac{1}{m} \sum_{k=1}^{m} S_k \tag{1}$$

where each scenario is evaluated through the involved actions

$$S_k(A_1^k, \ldots, A_l^k) = \frac{1}{l} \sum_{i=1}^{l} A_i^k \tag{2}$$

[7]https://docs.qameta.io/allure/



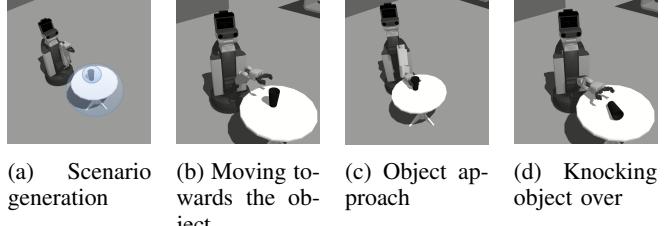| (a) Scenario generation | (b) Moving towards the object | (c) Object approach | (d) Knocking object over |

Fig. 5: Failed grasp execution from a coffee table

and each action is evaluated through its associated properties:

$$A_i^k(pp_1^i, \ldots, pp_n^i) = \frac{1}{n} \sum_{j=1}^{n} \mathbb{I}_{pp_j^i=1} \tag{3}$$

## IV. EXPERIMENTS

To verify our property-based testing approach, we test the Toyota HSR and use Gazebo [21] as a simulation environment since it is directly supported for the HSR.[8] We use ROS [22] as the underlying middleware, such that the tests are targeted at the components in our ROS-based domestic robotics stack.[9] ROS is thus used both for inter-component communication as well as for extracting information necessary for executing the property-based tests (such as object poses from Gazebo). We execute all tests on a laptop with an i5-6300HQ CPU at 2.30GHz and 8GB RAM running Ubuntu 16.04 and ROS Kinetic Kame. We present results for two use cases: testing an object grasping action and an object pick-and-place test. For both cases, we report the number of property successes and failures and include cases in which a test run could not be completed due to ROS communication issues[10] during the execution of an action (the latter are referred to as de-synced messages in the figures).

### A. Use Case 1: Object Grasping Action

Our first use case is that of verifying our *pick object* action. Each pick action test scenario follows four main steps, which are illustrated in Fig. 5 in the case of a failed grasp. During scenario generation, an object platform is spawned along with an object on it. For consistency, we performed this test with a small coffee table as an object platform and a glass as the object to be grasped. Once the test scenario is generated, the action is invoked and executed; this is followed by a verification of the properties in Table I. MoveIt[11] is used for trajectory planning and execution in simulation.

### B. Use Case 2: Complex Scenario

The second use case involves a complete object pick-and-place scenario, in which the robot navigates to a table, performs object detection on the table, picks an object, and

[8]The Gazebo setup for the HSR is available at https://github.com/hsr-project

[9]The central component of our domestic robotics stack can be found at https://github.com/b-it-bots/mas_domestic_robotics

[10]This particularly refers to action triggering messages that were not received by an action execution component, in which case the robot is waiting for a message indefinitely and thus remains stuck.

[11]https://moveit.ros.org

(a) Scenario generation

(b) Turning towards the table

(c) Table and object detection

(d) Cup grasping
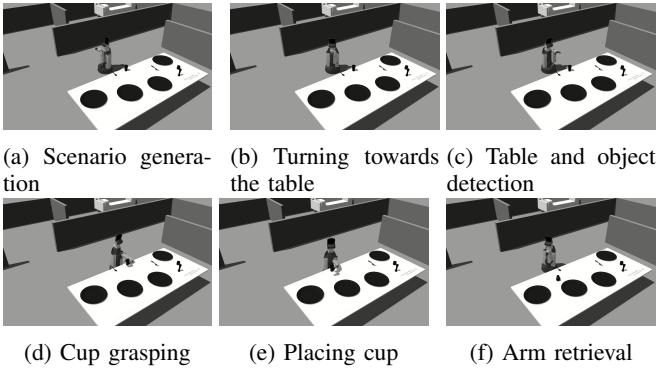
(e) Placing cup

(f) Arm retrieval

Fig. 6: A successful pick-and-place run. Both grasping and placing are performed on the same table.
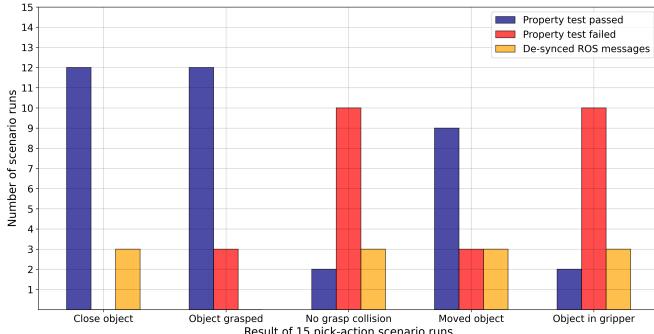


Fig. 7: Results of the pick action tests

places it back at a different location on the same table, as illustrated in Fig. 6. The objective of this use case is to validate the robot's behavior when performing a complete task that involves the actions *move to*, *perceive plane*, *pick object*, and *place object*, as the overall success in this task is determined by the sequential success of individual actions, namely task failures are likely to occur due to interdependencies between the actions.

### C. Results

Fig. 7 shows the results of 15 runs of the grasping action. As the results show, the action was fully successful only a few times. Most failures were caused by collisions with the object to be grasped while the approach trajectory was executed; as a result of these collisions, the object was displaced or knocked down and could not be successfully grasped. In some of the runs, the test could not be completed due to lost ROS messages between the components.

The results of 15 runs of the pick-and-place task are shown in Fig. 8. In this case, the majority of the tests failed due to lost communication messages, particularly during execution of the *pick* and *place* actions. Some of the tests failed while the robot was moving towards the object surface, which was caused by obstacles lying on the designated navigation goal; during execution of the *pick* and *place* actions, a collision was the most common failure cause, just as in the *pick* action test. For completeness, a detailed breakdown of the quantitative evaluation using Eq. 1-3 is provided in Table IV.
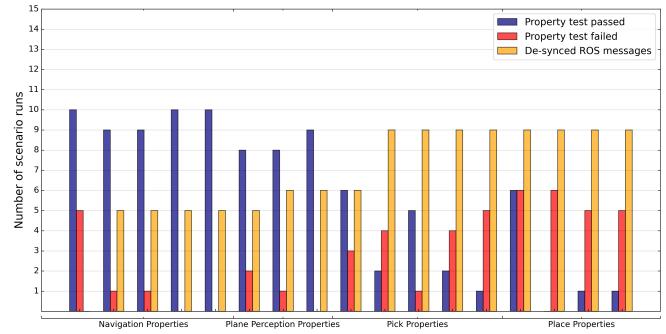


Fig. 8: Results of the pick-and-place scenario tests. Properties corresponding to different actions are separated into groups.

TABLE IV: Quantitative evaluation of 15 runs of the pick-and-place scenario. $A_1$-$A_4$ represent the actions *move to*, *perceive plane*, *pick object*, and *place object*, respectively.

| **Run** | $A_1^k$ | $A_2^k$ | $A_3^k$ | $A_4^k$ | $S_k$ |
|---|---|---|---|---|---|
| 1 | 1.0 | 1.0 | 0.4 | 0.25 | 0.67 |
| 2 | 1.0 | 0.67 | 0.0 | 0.0 | 0.42 |
| 3 | 0.6 | 0.33 | 0.6 | 0.25 | 0.45 |
| 4 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 5 | 1.0 | 1.0 | 0.0 | 0.0 | 0.5 |
| 6 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 7 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 8 | 1.0 | 1.0 | 0.4 | 0.25 | 0.67 |
| 9 | 1.0 | 1.0 | 0.4 | 0.25 | 0.67 |
| 10 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 11 | 1.0 | 1.0 | 0.4 | 0.25 | 0.67 |
| 12 | 1.0 | 1.0 | 1.0 | 0.75 | 0.94 |
| 13 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 14 | 1.0 | 1.0 | 0.0 | 0.0 | 0.5 |
| 15 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| | | | | $T$ | 0.37 |

It should be mentioned that the objective of the experiments is evaluating the effectiveness of our proposed approach rather than our components as such.[12] From this point of view, configurable property-based testing is able to consistently identify failed actions and generates useful information for finding the likely causes of execution failures.

## V. DISCUSSION AND CONCLUSIONS

In this paper, we proposed a framework for property-based testing of robots in simulation, with a particular focus on tabletop manipulation for a domestic robot. The framework utilises (i) a simulator that allows scenario randomisation and has information about the absolute state of the world, (ii) an ontology for encoding information about relevant action properties and parameters associated with those, and (iii) a configurable scenario generation component. The combination of these components allows testing robot actions, as well as action sequences, in diverse world configurations, which results in test reports that can be used to identify common failure patterns. In experiments with an object grasping action and a pick-and-place task, we showed that property-based testing in simulation is a promising approach that can

---

[12]Particularly since some of the components that we use on the real robot have not been adapted to the simulation.

potentially increase the transparency of service robots and provides insights about the reliability of tested components.

The work in this paper is only a preliminary step towards the use of property-based tests in simulation for regular robot testing. As can be seen from the experimental results, a considerable number of tests failed due to unsuccessful component communication through ROS; we have pinpointed this issue to limitations of the hardware used for running the tests, namely the resource-intensiveness of the simulation was negatively affecting the performance of the complete system, which means that more powerful hardware should be used for running tests so that they are more representative of the real system's behaviour. The Gazebo simulator, which we used for running the tests, was another frequent cause of failures; in particular, Gazebo crashes prevented several test runs from completing successfully, so it would be worthwhile to consider using alternative simulators instead [23]. Another limitation of our current approach is that some of the components that we use in simulation, such as the trajectory execution component, are modified versions of those used on the real robot due to minor discrepancies between the available interfaces on the real robot and in simulation; for the tests to correspond to the real system more closely, the interfaces in simulation and on the robot should be made as close as possible to each other. The realism of the tests is another potential concern for the proposed approach; this is particularly the case in the context of object interaction [24], although that would not be a significant problem if the simulated tests are only used as a preliminary step before running more critical tests on a real system and, depending on the tested policies, direct transfer to the real world may be possible in certain cases, as for instance shown in [25].

Future work will focus on testing execution policies that have been learned with real-world data, particularly those in [26], in diverse scenarios. Another aspect that should be addressed is that of improving the scenario generation and expanding the framework so that scenarios other than tabletop manipulation can be tested; this includes incorporating tests in which action sequences are generated by a task planner rather than being predefined during scenario generation. Adding additional scenario configuration parameters, particularly ones that control the physical properties of the objects used in the scenarios, would also contribute to the extensiveness and realism of the tests. Finally, it would also be useful to investigate post-processing of the generated test results so that the data can be utilised for simulation-based learning and policy improvement.

## REFERENCES

[1] T. Sotiropoulos, H. Waeselynck, J. Guiochet, and F. Ingrand, "Can robot navigation bugs be found in simulation? An exploratory study," in *Proc. IEEE Int. Conf. Software Quality, Reliability and Security (QRS)*, July 2017, pp. 150–159.

[2] A. Hentout, A. Mustapha, A. Maoudj, and I. Akli, "Key challenges and open issues of industrial collaborative robotics," in *RO-MAN Workshop on Human-Robot Interaction: From Service to Industry*, Aug. 2018.

[3] A. Mitrevski, A. Kuestenmacher, S. Thoduka, and P. G. Plöger, "Improving the reliability of service robots in the presence of external faults by learning action execution models," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2017, pp. 4256–4263.

[4] D. Araiza-Illan, A. G. Pipe, and K. Eder, "Model-based Test Generation for Robotic Software: Automata versus Belief-Desire-Intention Agents," *arXiv preprint arXiv:1609.08439*, 2016.

[5] F. George and B. Matt, "Property-Based Testing: A New Approach to Testing for Assurance," *ACM SIGSOFT Software Engineering Notes*, vol. 22, no. 4, pp. 74–80, 1997.

[6] T. Yamamoto, K. Terada, A. Ochiai, F. Saito, Y. Asahara, and K. Murase, "Development of Human Support Robot as the research platform of a domestic mobile manipulator," *ROBOMECH Journal*, vol. 6, pp. 1–15, 2019.

[7] P. Arora and R. Bhatia, "A Systematic Review of Agent-Based Test Case Generation for Regression Testing," *Arabian Journal for Science and Engineering*, vol. 43, pp. 1–24, Aug. 2017.

[8] A. Bihlmaier and H. Wörn, "Robot unit testing," in *Proc. 4th Int. Conf. on Simulation, Modeling, and Programming for Autonomous Robots (SIMPAR)*, vol. 8810, 2014, pp. 255–266.

[9] R. Majumdar, A. Mathur, M. Pirron, L. Stegner, and D. Zufferey, "Paracosm: A Test Framework for Autonomous Driving Simulations," in *Proc. 24th Int. Conf. Fundamental Approaches to Software Engineering (FASE)*, Feb. 2021, pp. 172–195.

[10] J. Park, M. Wen, Y. Sung, and K. Cho, "Multiple Event-Based Simulation Scenario Generation Approach for Autonomous Vehicle Smart Sensors and Devices," *Sensors*, vol. 19, no. 20, p. 4456, Oct. 2019.

[11] Y. Abeysirigoonawardena, F. Shkurti, and G. Dudek, "Generating Adversarial Driving Scenarios in High-Fidelity Simulators," in *Proc. Int. Conf. Robotics and Automation (ICRA)*, 2019, pp. 8271–8277.

[12] M. Koren and M. Kochenderfer, "Efficient Autonomy Validation in Simulation with Adaptive Stress Testing," in *Proc. IEEE Intelligent Transportation Systems Conference (ITSC)*, Oct. 2019, pp. 4178–4183.

[13] A. Santos, A. Cunha, and N. Macedo, "Property-based testing for the robot operating system," in *Proc. 9th ACM SIGSOFT Int. Workshop Automating Test Case Design, Selection, and Evaluation (A-TEST)*, 2018, pp. 56–62.

[14] V. Estivill-Castro, R. Hexel, and C. Lusty, "Continuous integration for testing full robotic behaviours in a GUI-stripped simulation," in *Proc. CEUR Workshop*, vol. 2245, 2018, pp. 453–464.

[15] L. Kunze, M. E. Dolha, E. Guzman, and M. Beetz, "Simulation-Based Temporal Projection of Everyday Robot Object Manipulation," in *Proc. 10th Int. Conf. Autonomous Agents and Multiagent Systems (AAMAS)*, 2011, pp. 107–114.

[16] L. Mösenlechner and M. Beetz, "Fast temporal projection using accurate physics-based geometric reasoning," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2013, pp. 1821–1827.

[17] J. Winkler, M. Tenorth, A. K. Bozcuoglu, and M. Beetz, "CRAMm - Memories for Robots Performing Everyday Manipulation Activities," *Advances in Cognitive Systems*, vol. 3, pp. 47–66, 2014.

[18] M. Wen, J. Park, and K. Cho, "A scenario generation pipeline for autonomous vehicle simulators," *Human-centric Computing and Information Sciences*, vol. 10, no. 1, June 2020.

[19] J. Arnold and R. Alexander, "Testing autonomous robot control software using procedural content generation," in *Proc. 32nd Int. Conf. Computer Safety, Reliability, and Security (SAFECOMP)*, vol. 8153, 2013, pp. 33–44.

[20] D. MacIver, Z. Hatfield-Dodds, and M. Contributors, "Hypothesis: A new approach to property-based testing," *Journal of Open Source Software*, vol. 4, no. 43, p. 1891, Nov. 2019.

[21] N. Koenig and A. Howard, "Design and use paradigms for Gazebo, an open-source multi-robot simulator," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems (IROS)*, 2004, pp. 2149–2154.

[22] M. Quigley *et al.*, "ROS: an open-source Robot Operating System," in *ICRA Workshop on Open Source Robotics*, May 2009.

[23] J. Collins, S. Chand, A. Vanderkop, and D. Howard, "A Review of Physics Simulators for Robotic Applications," *IEEE Access*, vol. 9, pp. 51 416–51 431, 2021.

[24] J. Collins, D. Howard, and J. Leitner, "Quantifying the Reality Gap in Robotic Manipulation Tasks," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2019, pp. 6706–6712.

[25] L. Hermann *et al.*, "Adaptive Curriculum Generation from Demonstrations for Sim-to-Real Visuomotor Control," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2020, pp. 6498–6505.

[26] A. Mitrevski, P. G. Plöger, and G. Lakemeyer, "Representation and Experience-Based Learning of Explainable Models for Robot Action Execution," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems (IROS)*, 2020, pp. 5641–5647.