

# Fidelity of Cloud Emulators: The Imitation Game of Testing Cloud-based Software

Anna Mazhar<sup>§</sup>, Saad Sher Alam, William Zheng, Yinfang Chen, Suman Nath<sup>†</sup>, Tianyin Xu

<sup>§</sup>Cornell University, Ithaca, NY, USA

University of Illinois Urbana-Champaign, Urbana, IL, USA

<sup>†</sup>Microsoft Research, Redmond, WA, USA

**Abstract**—Modern software projects have been increasingly using cloud services as important components. The cloud-based programming practice greatly simplifies software development by harvesting cloud benefits (e.g., high availability and elasticity). However, it imposes new challenges for software testing and analysis, due to opaqueness of cloud backends and monetary cost of invoking cloud services for continuous integration and deployment. As a result, cloud emulators are developed for offline development and testing, before online testing and deployment.

This paper presents a systematic analysis of cloud emulators from the perspective of cloud-based software testing. Our goal is to (1) understand the discrepancies introduced by cloud emulation with regard to software quality assurance and deployment safety and (2) address inevitable gaps between emulated and real cloud services. The analysis results are concerning. Among 255 APIs of five cloud services from Azure and Amazon Web Services (AWS), we detected discrepant behavior between the emulated and real services in 94 (37%) of the APIs. These discrepancies lead to inconsistent testing results, threatening deployment safety, introducing false alarms, and creating debuggability issues. The root causes are diverse, including accidental implementation defects and essential emulation challenges. We discuss potential solutions and develop a practical mitigation technique to address discrepancies of cloud emulators for software testing.

## I. INTRODUCTION

Modern software projects have been increasingly using cloud services as important components for storage, database, data processing, etc. Such cloud-based programming practice greatly simplifies software development by harvesting cloud benefits (e.g., high availability and elasticity) and software deployment by reducing the cost of purchasing and managing large-scale systems and infrastructures. Today, all major cloud providers offer various services to support cloud-based software and these cloud services are widely used [1]–[3].

Despite its benefits, cloud-based programming imposes new challenges for software testing and analysis due to opaqueness of cloud backends and monetary cost of invoking cloud services during continuous integration and deployment (CI/CD). First, unlike other types of dependencies like libraries, which are linked as a part of the software program, cloud services are *external* to cloud-based software (invoked via REST API calls), and their backend implementations are opaque. It is hard to reason about the correctness of cloud-based software independently, especially its end-to-end behavior. For example, regressions of cloud backend implementations [4] can directly affect dependent software that invokes corresponding APIs.

Second, testing cloud-based software with real services can be costly, especially with CI/CD. Cloud services charge users based on the number of API invocations, storage, and additional features like transaction support [5], [6]. So, extensive testing on the cloud is expensive. For example, the test suite of Orleans issues 120K+ Azure API calls. Under CI/CD, tests are continuously invoked [7]–[9]. We expect even higher costs in the near future as cloud services are increasingly adopted by software projects and new tests are being added.

*Cloud emulators* are developed for cloud-based software development and testing before online testing and deployment. Nine out of ten projects we studied (§III) use emulators for CI tests. We are also informed by a major cloud service provider that emulators are widely used by customers who use their service APIs. A cloud emulator offers local simulation of large, complex cloud services. For example, a fault-tolerant, persistent key-value storage service can be emulated by a centralized, in-memory hash table [10]. Cloud emulators enable developers to conduct prompt, cost-efficient offline testing and debugging [11]. They are transparent to software under test—using emulators requires no code change but a simple setup to connect to emulated services. Cloud emulators are typically developed or supported by cloud service providers. For example, Microsoft provides emulators for Azure services, e.g., Azurite [12] for Azure Storage Services [13].

Ideally, emulators should behave the same as real cloud services so that software quality assurance, like testing, can rely on emulators. However, it is prohibitively difficult for emulators to achieve perfect fidelity (considering the complexity, scale, and distributed nature of cloud services). In practice, emulators implement specifications of cloud service APIs (§II). However, as shown in our study (§V-A), specifications of today’s cloud services and their APIs are often incomplete and limited. Without formal enforcement of emulator compliance with real cloud service, it is unclear how much fidelity today’s emulators could realize. We use the term *discrepancies* to refer to emulator behavior that deviates from specified behavior of cloud services. We observed that discrepancies are constantly reported to affect testing of cloud-based software [14]–[17].

In this paper, we analyze the discrepancies between cloud emulators and cloud services to understand the fidelity of cloud emulation in practice and its impacts on cloud-based software quality assurance and developer experience. Specifically, we apply differential testing against two widely used cloud em-

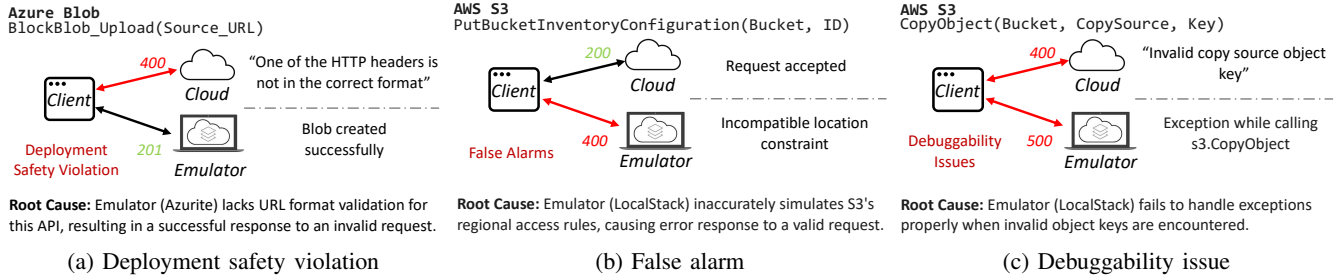


Fig. 1: Implications of discrepancies between cloud emulators and cloud services with regards to software testing.

ulators, Azurite [12] for Azure Storage services (including Blob, Table, and Queue) and LocalStack [18] for Amazon Web Services (including S3 and DynamoDB). We record discrepant behavior between the cloud emulator and the cloud services, and analyze the root cause of each discrepancy. We focus on basic functional correctness, instead of performance or fault tolerance (e.g., data consistency and crash consistency) which are beyond the expectation of local emulation.

Our analysis results are concerning. Among 255 APIs of five cloud services from Azure and Amazon Web Services (AWS), we detected discrepant behavior between the emulated and real services in 94 (37%) of the APIs. These discrepancies have profound implications on deployment safety and developer experience: (1) code that passes tests with emulators may fail in production when cloud services are used; (2) test failures with emulators can be false alarms; and (3) debugging with emulators can be hard due to discrepant feedback (e.g., error code and messages). Figure 1 shows three examples we discovered in our analysis. We further analyze ten open-source cloud-based software projects; five of them are affected by discrepancies—some of their tests have inconsistent results when running on the cloud emulator versus the cloud services. In one project (Durabletask [19]), 78% of the tests are affected.

The root causes of discrepancies are diverse but can be categorized into (1) incompleteness of existing specifications, (2) unspecified behavior, and (3) implementation defects (such as bugs and missing features). While these root causes reflect essential software engineering challenges, we believe that many discrepancies could be addressed by more comprehensive testing and more systematic specification. We discuss potential solutions and mitigations, ranging from practical formal methods to new system-level support (§VI).

We explore hybrid cloud-emulator testing as a short-term mitigation and develop a simple tool named ET to selectively run tests on emulators versus cloud services, based on whether the test invokes discrepant APIs (§VII). ET offers different policies depending on whether discrepant API information is known as apriori or being done via in-situ analysis. Through ET, we show that hybrid testing yields considerable cost savings compared to running all tests with cloud services.

The paper makes the following main contributions:

- A discussion on the challenges of testing cloud-based software and the discrepancies introduced by emulators.
- A systematic analysis of discrepancies between cloud ser-

vices and their emulators, including their characteristics, root causes, and impacts on software testing.

- A discussion on solutions to discrepancies and a mitigation tool to selectively run tests with emulators.
- We reported bugs that caused the discrepancies; so far, six have been confirmed and five have been fixed.
- Research artifact: <https://github.com/xlab-uiuc/cloudtest>.

## II. BACKGROUND

### A. Cloud Services and Their APIs

Modern cloud services are programmatically accessed via REST APIs (defined as follows) on top of HTTP(S).

**Definition 1 (REST API).** In this paper, we defined a REST API as an HTTP method plus the resource. For example, “GET /blog/posts/{id}” and “PUT /blog/posts/{id}” are considered two unique REST APIs of a web blogpost service.<sup>1</sup>

A service typically exposes several tens of REST APIs. For example, AWS S3 exposes 97 REST APIs [22], and Azure Blob service exposes 72 REST APIs [23]. To ease developer programming, cloud services provide Software Development Kits (SDKs) with high-level, language-specific library APIs. Typically, SDK APIs wrap raw REST APIs, and for a service, more SDK APIs are built on top of the REST APIs. Most existing cloud-based applications invoke SDK APIs to interact with the cloud services instead of calling raw REST APIs.

1) *API Specification:* The REST APIs are commonly described using specification languages such as OpenAPI Specification [24]. The specification describes the API version, request URI, content type, input parameter, output format, error code and messages, etc. The API specifications are used by cloud emulators (§II-B) to develop emulated APIs.

We find that the API specifications are often incomplete. For example, the OpenAPI specification of Azure Blob services only specifies value constraints (including data types) for 63 (59%) of 107 parameters across all Azure Blob APIs. Besides, all the specifications are on data type, value range, and default value with no behavior semantics (e.g., async or not).

SDKs often include additional checks on parameter values of API calls over the API specifications—values that satisfy the API specification could be rejected by the SDK checks.

<sup>1</sup>We follow the REST API definition used in Azurite [20]. This definition can be inconsistent with other definitions, e.g., an API in our definition is referred to as an “API operation” in [21]. We choose this definition because it resembles SDK APIs faced by developers (e.g., “read a blog post” and “create a new blog post” are corresponding to two different SDK APIs).

2) *Pricing*: Cloud services are expensive. Despite different pricing models of cloud services, pricing typically depends on the amount of data to be stored and the cost of operations. Take Azure Blob service as an example. The price for 100TB/month ranges from \$91–\$1,545, depending on the access tiers [5]. Azure Blob service then charges for read, write, iterative-read, and iterative-write operations separately [25]. For example, the price for write operations varies from \$0.0228–\$0.13 per 10,000 writes, depending on the tier. Other features, such as redundancy [26]–[28], further increase the cost.

With the current pricing model, testing cloud-based software incurs non-trivial monetary costs. To demonstrate the cost, we run the tests of Orleans (a cloud-based software project) for Azure Storage services with standard configuration. Orleans has 189 tests that issue 120K+ Azure API calls over 23 unique APIs. We run these tests 500 times, which costs \$74.5 US dollars (we expect 500 times to be a reasonable time in CI/CD of large software projects [29], [30]).

### B. Cloud Emulator

To reduce cost and get prompt feedback, emulators are developed to assist developers in offline development and testing. Emulators can also be used for debugging production problems. Emulators run as local daemons that simulate cloud services. Cloud-based software programs transparently interact with the emulator in the same way they interact with cloud services. Using an emulator only needs a simple configuration that switches the connection from a cloud handle to localhost listened to by the emulator; no code change is needed.

Most cloud services provide developers with official emulators. For example, Microsoft provides emulators for Azure Storage and CosmosDB, and AWS provides emulators for DynamoDB and Step Functions. Moreover, third-party emulators are developed. One successful example is LocalStack [18], which emulates many AWS services such as S3 and DynamoDB. Compared with official emulators, LocalStack provides a more usable integrated development environment [31]. Our study deliberately selects an official emulator (Azurite) and a third-party emulator (LocalStack).

For compliance with the target cloud services, cloud emulators are commonly built on top of API specifications. For example, Azurite uses AutoRest [32] to generate stub code from the OpenAPI specification of Azure Storage services [33]. LocalStack employs weekly GitHub Action Checks to detect any changes of the API specifications of AWS [34].

### C. Emulation versus Mocking

Mocking is a common practice used in *unit tests* to simulate dependencies of code under test [35], [36]. Unlike emulators, mocked objects are not required to rigorously satisfy API specifications, because mock uses “behavior verification [36].” For this reason, mock cannot help with state verification.

Emulation is fundamentally different from mocking. An emulator is expected to conform to API specifications of the cloud services (for proprietary services, API specifications are the contract). The emulator is designed as a drop-in replacement

TABLE I: Emulators and cloud services studied in this paper (only AWS services studied in this paper are listed).

| Emulator   | Service            | LOC    | #Commits | Developer   |
|------------|--------------------|--------|----------|-------------|
| Azurite    | Blob, Queue, Table | 2,591K | 1,034    | Official    |
| LocalStack | S3, DynamoDB       | 449K   | 5,527    | Third-party |

for the actual service so developers can move from testing to deployment by simply changing a connection string. Emulators are closer to “fakes” [35], [36] than mocks. An emulator must maintain states so producer-consumer dependencies are expected to be the same. The emulator does not have to behave exactly the same or maintain the same internal states as the actual service beyond API specifications.

Without high-fidelity emulators, cloud-based software developers can only run tests with the actual services at certain time points during CI/CD. However, this creates a difficult trade-off between cost and effectiveness. Running tests against actual services frequently may incur high costs. Infrequent testing with actual services leads to big bundles of commits, affecting CI/CD effectiveness and reducing developer experience. Therefore, emulator fidelity is important.

## III. METHODOLOGY

We use differential testing to discover discrepancies between cloud emulators and real cloud services. Basically, we issue the same REST API calls to the emulated service and the cloud service independently and check the resulting behavior, including the return values, error codes or messages (if any), and states of key data objects such as blobs and containers. Any inconsistent behavior indicates a discrepancy.

**Studied emulators.** We select two cloud service providers with the highest market share, Microsoft Azure and Amazon Web Services (AWS) [37]. For these two providers, we choose to study the most commonly used emulators: Azurite [12] and LocalStack [18]; they represent state-of-the-art. Moreover, Azurite represents the official emulators provided by cloud service providers, while LocalStack represents third-party emulators developed by companies of cloud-based integrated programming environments. Importantly, both emulators are open-sourced, which enables us to debug discovered discrepancies. Table I lists the information of the two emulators.

**Studied services.** For the two emulators, we select five widely used cloud services: Blob, Queue, and Table services from the Azure Storage services and S3 and DynamoDB from AWS. Azurite only supports Azure Storage services (Blob, Queue, and Table); for LocalStack, we pick DynamoDB and S3 as popular and commonly used AWS services.

**Test workloads.** We use two complementary test workloads. First, we leverage API fuzzing to generate sequences of REST API calls. Each API call sequence is a test workload and the workloads collectively cover all the REST APIs provided by the target cloud services. The API fuzzing workloads help us understand the discrepancies in each REST API and characterize a broad range of APIs.

The fuzzing is done against SDK APIs, not raw REST APIs. We in fact started from REST API fuzzing using RESTler [21]. However, we found that certain discrepancies are not possible if the software under test uses SDKs which have additional checks (§II). In practice, developers do not commonly craft REST API calls directly but mostly call SDK APIs (§II-A). Since our goal is to understand discrepancies in the context of software development, rather than security analysis [21], [38], we choose to fuzz SDK APIs. Basically, we focus on analyzing discrepancies faced by cloud-based software developers.

We also use the test suites of existing cloud-based software projects as the test workloads. Many tests invoke cloud service APIs. These tests help understand the impact of discrepancies on testing real-world software projects, which is complementary to fuzzing from the API perspective.

**Fuzzing SDK APIs.** We implemented a grammar-aware API fuzzer to generate diverse SDK API calls as test workloads. We start from default or predefined parameter values for each SDK API and the fuzzer mutates parameter values based on value constraints defined in OpenAPI specifications of REST APIs (the “grammar”). To do so, we establish the mapping from the parameters of REST APIs to those of the corresponding SDK APIs; the mapping process is straightforward because SDK APIs are mostly wrappers over raw REST APIs. The grammar-based mutation ensures that generated SDK API calls are mostly valid and can reach emulated or real cloud services. Our fuzzer implements the fuzzing approach of RESTler [21]: (1) inferring producer-consumer dependencies among request types (e.g., “API *Y* should be called after API *X*” because *Y* takes as an input a resource-ID produced by *X*) and (2) taking dynamic feedback from responses during testing (e.g., learning that “a API *Y* called after a sequence  $X \rightarrow Y$  is refused” and avoiding this combination in the future).

We monitor the response of each API call. If inconsistent responses are returned by the emulator and the cloud services (including both HTTP response status code like 200 and 404, as well as error code and message if the response returns an error), we capture and record the discrepancy and abort the test. Otherwise, we progress to the next API call in the generated sequence. We also check the key data objects before and after the API calls (e.g., the number of blobs for Azure Blob Services) to capture discrepancies with no immediately observable manifestation, such as resource leaks (§IV-B). Those checks are service-specific.

**Using existing tests.** To understand the impact of discrepancies on real-world software projects, we perform differential testing using test suites of existing projects. We select ten open-source projects (Table II) that use the studied cloud services. The ten projects are selected because they are mature and widely used (based on the total number of commits and star counts), developed by reliable sources such as companies like Microsoft (Orleans, DurableTask), NuGet (Insights), and PetaBridge (Alpakka), and are actively maintained and use recent versions of the studied cloud service APIs.

In our study, we select tests that interact with the cloud

TABLE II: Cloud-based software projects. “#Tests” refers to tests that invoke cloud services; “#APIs” refers to *unique* APIs.

| Project            | Services            | LOC    | #Tests | #APIs |
|--------------------|---------------------|--------|--------|-------|
| Alpakka            | Queue               | 22.5K  | 9      | 6     |
| AttachmentPlugin   | Blob                | 1.9K   | 23     | 7     |
| DurableTask        | Blob, Queue, Table  | 59.0K  | 101    | 30    |
| IdentityAzureTable | Table               | 85.7K  | 51     | 6     |
| Insights           | Blob, Queue, Table  | 144.8K | 171    | 20    |
| IronPigeon         | Blob                | 37.8K  | 7      | 8     |
| Orleans            | All services but S3 | 204.8k | 247    | 35    |
| ServiceStack       | DynamoDB, S3        | 756.2K | 187    | 15    |
| Sleet              | Blob, S3            | 21.2K  | 22     | 21    |
| Streamstone        | Table               | 4.6K   | 75     | 7     |

services. The selection is done by monitoring the HTTP traffic of each test in a reference run using the emulator. We check whether a test outputs inconsistent results when running with emulators versus cloud services. Table II shows the number of tests that invoke cloud services and the number of unique APIs invoked by the test suite of each project.

**Discrepancy analysis process.** We inspect every observed discrepancy during the aforementioned testing. A discrepancy is recorded if it manifests via inconsistent test results—the results of a test are different when running with the emulator versus the real services. For each discrepancy, we verify it by deterministically reproducing its manifestation and impact, and debug it to localize the root cause in emulator source code. The process helped us minimize human errors during the analysis and subjectiveness in the interpretation and categorization.

#### IV. DISCREPANCY CHARACTERISTICS

##### A. Prevalence of Discrepancies

Our analysis shows that discrepancies are prevalent in the two cloud emulators (Azurite and LocalStack). The five cloud services we studied expose a total of 255 APIs. Among these 255 APIs, our API fuzzer (§III) discovered discrepancies in 94 (37%). Table III shows the number of discrepant APIs of each service. We define a *discrepant API* as follows:

**Definition 2 (Discrepant API).** An API is discrepant if it can expose inconsistent behavior, in the scope of its specification, when invoked on the emulator versus on the real cloud service.

Both Azurite and LocalStack have a considerable percentage of discrepant APIs among all the APIs they support and across the services, showing that discrepancies are not specific to one emulator implementation or specific to APIs of a particular service. Rather, emulator fidelity is a common challenge.

These discrepancies have different implications as exemplified in Figure 1, including (1) deployment safety violations (1a), (2) false alarms (1b), and (3) debuggability issues (1c). Table IV categorizes the implications of the total 98 discrepancies discovered in the 94 discrepant APIs (one discrepancy can have different implications). The results show diverse implications of these discrepancies. We measure their impacts on real-world test cases in §IV-C.

Surprisingly, we find that 37 of the 94 discrepant APIs are certified by the emulators and considered “fully supported.”

TABLE III: Discrepant APIs with respect to the cloud services

| Services     | Emulator   | Total APIs | Discrepant APIs |
|--------------|------------|------------|-----------------|
| Azure Blob   | Azurite    | 72         | 31 (43%)        |
| Azure Table  | Azurite    | 15         | 1 (7%)          |
| Azure Queue  | Azurite    | 18         | 2 (11%)         |
| AWS S3       | LocalStack | 97         | 33 (34%)        |
| AWS DynamoDB | LocalStack | 53         | 27 (51%)        |
| <b>Total</b> |            | 255        | 94 (37%)        |

TABLE IV: Impacts of discrepancies across emulators.

| Impact               | Azurite/Azure | LocalStack/AWS | Total |
|----------------------|---------------|----------------|-------|
| Deployment safety    | 13            | 22             | 35    |
| False alarms         | 12            | 33             | 45    |
| Debuggability issues | 9             | 9              | 18    |
| <b>Total</b>         | 34            | 64             | 98    |

LocalStack adopts five methods to certify emulated APIs [39], [40], including both internal and external integration tests (e.g., snapshot tests [41]). Despite extensive efforts, 22 (out of 60) discrepant APIs from LocalStack are certified by all five testing methods, while 39 (out of 60) discrepant APIs are certified with at least one test method. Similarly, 15 (out of 34) discrepant APIs from Azurite are certified to be fully supported by Azurite [20]. The results show the challenges faced by existing testing-based practices in detecting discrepancies.

**Finding 1.** *Discrepancies between modern cloud emulators and real cloud services are prevalent (discovered in 37% of APIs on average and even in certified APIs). The implications of discrepancies are diverse, including unsafe deployment, false alarms, and debuggability issues.*

### B. Discrepancy Manifestations

A notable observation is that discrepancies are manifested through not only inconsistent responses to the API calls, but also inconsistent remote, cloud-side states. The latter creates significant challenges to observe and understand discrepancies, especially with short-running test cases. We implemented domain-specific checks to compare the remote states maintained by the emulators and the corresponding cloud services (§III). For example, we check the states of each container (maintained by the emulators and cloud services) before and after each container-related API call.

Seven discrepancies have the same response to API calls but create inconsistent remote states. For example, when invoking an Azure Blob API, “Container\_Restore [42]”, to recover an early deleted container, both Azurite and the Azure Blob service return the same response; however, the Blob service faithfully restores the deleted container, while Azurite creates a new empty container. Such discrepancies may not be easy to capture without fine-grained checks.

We also find 44 discrepancies that cause inconsistent responses and inconsistent remote states. For example, when calling a Blob API “BlockBlob\_StageBlockFromURL” [43] with an invalid URL, Azurite succeeds by creating a new blob, while the Blob service fails with `InvalidHeaderValue`.

TABLE V: Software tests that are affected by discrepancies.

| Project      | # Discrepant Results | # Discrepant Tests | Impact      |              |
|--------------|----------------------|--------------------|-------------|--------------|
|              |                      |                    | Safety Vio. | False Alarms |
| Alpakka      | 9 (100%)             | 9                  | 9           | 0            |
| DurableTask  | 79 (78%)             | 101                | 79          | 0            |
| Orleans      | 8 (9%)               | 82                 | 5           | 3            |
| ServiceStack | 3 (2%)               | 72                 | 2           | 1            |
| Streamstone  | 1 (1%)               | 75                 | 1           | 0            |
| <b>Total</b> | 100                  | 339                | 96          | 4            |

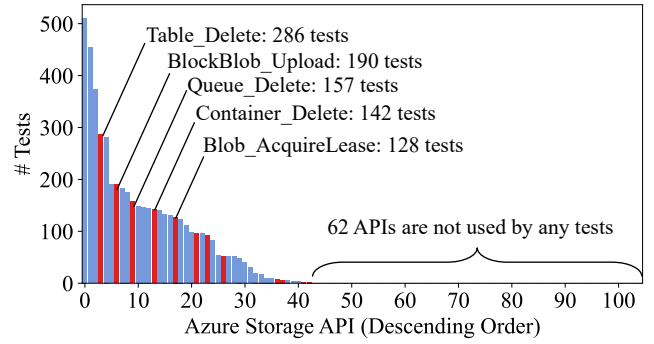


Fig. 2: Popularity of Azure Storage APIs, measured by the number of tests that use an API across all the studied projects.

**Finding 2.** *Discrepancies that only manifested via inconsistent remote states are hard to observe; fine-grained state checks are needed to capture those silent discrepancies.*

### C. Impact on Real-world Tests

We measure the impacts of discrepancies on real-world test suites of cloud-based software projects (Table II). The impact is reflected by inconsistent test results when running the same test with the emulator versus the cloud service. We define discrepant tests as follows:

**Definition 3 (Discrepant test).** A test is a discrepant test if it invokes any discrepant APIs. Note that a discrepant test may or may not output discrepant test results, depending on whether API call statements are executed and whether discrepancy-inducing parameters are used during the test execution.

Among the ten projects we evaluated (Table II), we discovered discrepant test results in 50% of them (five projects), as shown in Table V. Different projects are affected at different levels, ranging from 1% to 100% of tests that invoke cloud services. The variation is attributed to the usage characteristics of the cloud service APIs. Specifically, though we discovered a large number of discrepant APIs (§IV-A), not all these APIs are equally invoked by the test cases. Figure 2 depicts the popularity of all Azure Storage APIs (105 in total) invoked by all the tests of the studied projects, where discrepant APIs are marked in red. Popularity is measured by the number of tests that use the API. Among the 105 APIs, only 43 of them are invoked by at least one test. Only 12 APIs (out of 43) involved in the tests are discrepant (while 34 discrepant APIs in total are discovered for Azure; see Table III).

Note that the number of discrepant tests is much larger than the number of discrepant test results manifested (Table III).

The reason is that many discrepant tests are only manifested when certain parameter values are used.

The results have two important implications. First, addressing discrepancies can leverage API usage characteristics in the field to prioritize widely used APIs. Oftentimes, fixing discrepancies of a few APIs can eliminate a large number of discrepant tests or test results. We take all the tests using Azure services as the example: by resolving the top five discrepant Azure APIs in Figure 2, discrepant tests drop by 63% (from 267 to 99) and discrepant results drop by 10% (from 89 to 80). The small drop in discrepant results is caused by tests in DurableTask utilizing multiple discrepant APIs. If the top seven discrepant Azure APIs are resolved, 75 out of the 79 discrepant results caused by DurableTask will be eliminated.

Second, fine-grained, parameter-level analysis can further capture discrepancies. Although our analysis stays at the API level instead of parameters, we build on these implications when designing mitigation solutions (§VII).

**Finding 3.** *Five out of ten studied software projects reveal discrepant test results caused by discrepancies between emulators and cloud services. Those discrepant tests are caused by a small set of discrepant APIs. Not all discrepant APIs manifest during testing if triggering parameters are not used.*

We further categorize the implications of discrepant tests into (1) *deployment safety violations* (1a), and (2) *false alarms* (1b), as broken down in Table V. Debuggability issues are not applicable here as the tests all pass in the default setup.

The majority of test discrepancies would lead to deployment safety violations—the test that passes with the emulator would fail when running with the cloud service (i.e., passing the test provides no safety guarantee on the cloud). For example, a test `CreateTaskHub` in DurableTask uses the Azure Blob API, `Container_Create`, to re-create a previously deleted blob container. This test fails when running with the cloud service due to `DurableTaskStorageException`—“*the specified container is being deleted; try operation later*,” because container deletion is asynchronous and provides no guarantee for the time to finish. However, this test always passes when running with the Azurite emulator, as Azurite always deletes the container synchronously before the API returns.

False alarms are relatively less common than deployment safety issues (Table V). Two (out of four) false alarms are caused by brittle assertions on the error messages returned by the API calls (which are discrepant between the emulator and the cloud service). Such discrepancies can be addressed by enforcing the consistencies of the error messages. One false alarm is caused by a flaky test [44], [45]; the non-deterministic flaky behavior only manifests when running with the emulator, not with the cloud service, due to order differences caused by discrepant timing of API calls. We fixed the flaky tests by adding `await` to enforce the order. The last false alarm is caused by resource discrepancy—the `stress-test` in Orleans exhausted the socket limit of `LocalStack` (which passes with the cloud service). Such resource discrepancies are essential,

TABLE VI: Root causes of observed discrepancies.

| Service      | Incomplete Spec. | Unspecified | Defects in Impl. |
|--------------|------------------|-------------|------------------|
| Azure Blob   | 18 (58.1%)       | 1 (3.2%)    | 12 (38.7%)       |
| Azure Queue  | 1 (50.0%)        | 1 (50.0%)   | 0 (0.0%)         |
| Azure Table  | 0 (0.0%)         | 1 (100.0%)  | 0 (0.0%)         |
| AWS S3       | 11 (29.7%)       | 14 (37.8%)  | 12 (32.4%)       |
| AWS DynamoDB | 2 (7.4%)         | 4 (14.8%)   | 21 (77.8%)       |

and stress tests should not use emulators in the first place.

**Finding 4.** *Deployment safety violations are the major implications of discrepant tests, while false alarms also appear in testing results. Tests of cloud-based software projects need to carefully decide to run on emulators versus cloud services.*

## V. ROOT CAUSE ANALYSIS

We discuss the discrepancies from the specification perspective. Conceptually, both the emulator and the cloud services are implementations of the API specification (in practice, emulators implement the API specifications defined by the cloud services). So, discrepancies are the result of defects in either specification or implementation. Based on the existing API specifications (§II), we categorize the discrepancies into: (1) incomplete specification, (2) unspecified behavior that is not considered in existing specifications, and (3) implementation defects in the emulators or the cloud services. Table VI shows the three categories of discrepancy root causes.

During the project, we detected ten bugs in the two studied emulators, of which six have been confirmed (and five fixed). We also detected two bugs in the cloud backend implementations, which have been reported to the cloud service providers.

### A. Incomplete Specifications

As discussed in §II, existing cloud service API specification focuses on parameter value constraints and error codes and messages, from which emulators automatically generate stub code that adheres to the specifications (e.g., using `AutoRest` [32]). However, we still find that a significant percentage of discrepancies are caused by inconsistent validity checks of parameter values as well as inconsistent error code and messages. The reason is incomplete specifications. Table VI shows that incomplete specifications can cause up to 58.1% of the discovered discrepancies in a service.

1) *Parameter value constraint*: Ideally, the API specification should define *all* the value constraints of every input parameter. In reality, API specifications are deficient. We observe discrepant value constraint checks in twelve out of 34 discrepant Azure Storage APIs and ten out of 33 AWS S3 APIs. We observe no such discrepancy in DynamoDB.

Figure 3a shows such an example from Azure Blob, where the value of the `x-ms-proposed-lease-id` parameter of the `Blob_ChangeLease` API should be in the GUID format [46]. The Blob service implements a format check, while Azurite does not. As a result, an invalid API call of `Blob_ChangeLease` will be returned successfully by the emulator but rejected by



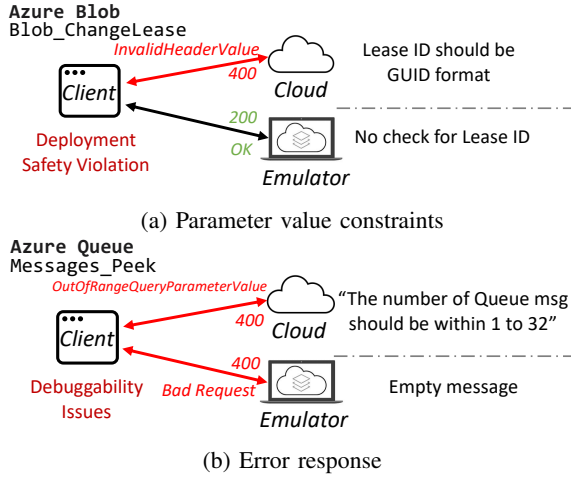


Fig. 3: Discrepancies caused by deficient specifications

the cloud. However, such value constraint is not specified in the OpenAPI specification of Azure Blob services.

In another common discrepancy case across Azure and AWS, cloud service APIs require authorization to private resources or sensitive operations (e.g., security configuration like `PutBucketAcl`). In its absence, these requests are denied by the cloud service. However, our results revealed that emulators often overlook this constraint, accepting such requests with a 200 OK response, resulting in 8% of discrepancies. We find that the requirement of authentication is commonly included in the text descriptions, which is not machine-checkable. Although this is not a defect as far as OpenAPI is concerned, it is not enforced by auto-generated stub code. From a codegen/machine-checkability perspective, the specification is incomplete. Our experience of examining Azure and AWS OpenAPI specifications shows that text-based API descriptions often includes constraints that are not machine-checkable.

2) *Error response*: We also find that specifications can be incomplete in the expected error code and messages and fail to associate them with the APIs, leaving emulator developers to interpret discrepant error messages. Figure 3b shows such an example. When a request is made with an out-of-range value for the `numofmessages` parameter, the Azure Queue service provides a detailed message pinpointing the error. In contrast, Azurite only responds with a “Bad Request” error code, offering no specific guidance and impeding debuggability.

Discrepant error responses were particularly prominent in Azure Storage APIs, accounting for 21% (7 out of 34) of the total discrepancies. When we examined the Azure API specifications, we found that the error codes were not associated with the APIs but were defined in a separate list. Differently, we found that AWS specifications have a more structured approach to error code definitions, which were also part of the related API definitions. The latter directly translates to the emulator code. In DynamoDB API specifications, we found structured definitions of 31 unique error codes, including their error messages, exception flags, and documentation. Hence, discrepant error responses are rare in DynamoDB and S3.

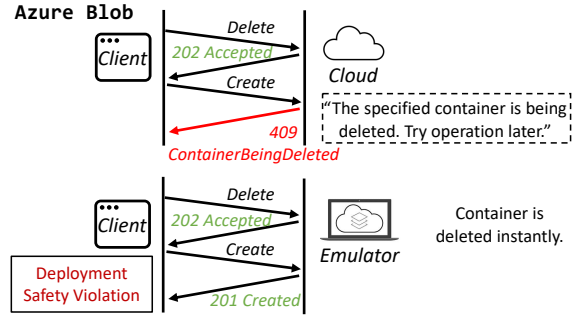


Fig. 4: Discrepancies caused by unspecified behavior

Techniques for generating and enforcing machine-checkable API specifications may potentially close the gaps. Recent work demonstrates the potential of synthesizing formal specifications from text [47]–[49]. Techniques that infer specifications from code [50]–[52] can help differential analysis of API specifications from emulator and service code to find discrepancies. Certainly, API is a form of abstraction. So, overly detailed specifications that describe internal behavior could be considered leaky abstractions. A key challenge is to define the right level of specifications as an effective API contract.

**Finding 5.** *The completeness of machine-checkable specifications is still a fundamental challenge, even for simple specifications such as parameter value constraints and error code. Without an effective way towards comprehensive specifications, we expect such discrepancies to remain prevalent.*

## B. Unspecified Behavior

A few discrepancies were caused by API behavior out of the scope of the existing specification and thus is unspecified. We observed two patterns of unspecified-behavior discrepancies.

We mentioned the first pattern in §IV-C—whether an API is synchronous or asynchronous. For example, Azure Blob’s API `Container_Delete`, which deletes container resources in cloud services, is an asynchronous API. For efficiency consideration, the deletion is not guaranteed to finish before the API returns. Instead, the time to finish the deletion depends on the amount of resources to be deleted. Conversely, emulators always finish deletions before returning the API calls. Figure 4 depicts such discrepancies. The result is that API sequences involving creating a container, deleting it, and then attempting to recreate it with the same name yielded different results: the cloud service returned 409 `ContainerBeingDeleted`, while the emulator allowed immediate container recreation with 201 `Created`. This pattern also appears in sequences following a deletion API call: the emulator would return a 404 `Not Found` after deletion, while the cloud, busy doing the deletion, would non-deterministically (depending on timing) issue a success response or a 409 `Conflict` message, “*The specified container is being deleted. Try operation later.*”

The second pattern is unspecified API behavior on null references (e.g., non-existent objects). For example, when using LocalStack, the emulated S3 APIs for fetching bucket configu-

ration (e.g., `GetBucketMetricsConfiguration`) or policy (e.g., `GetBucketPolicyStatus`) would return a 200 OK response with an empty policy configuration in the response, when configurations were never set. In contrast, the real S3 APIs respond with a 404 error, suggesting that the configuration was not found. A similar example is APIs responsible for deleting configurations (e.g., `DeleteBucketMetricsConfiguration`) or object tags (e.g., `DeleteBucketPolicy`). If a configuration was not created, the emulator responded with a 204 success upon deletion, while the cloud service returned a 404 error. For S3 on LocalStack, eleven discrepancies were caused by such a case. Such undefined behavior resembles null pointers as a common source of undefined behavior [53].

**Finding 6.** *Two patterns of undefined behavior contribute to discrepancies between emulators and cloud services: (1) the synchrony of the API and (2) null references. Such behavior is currently not considered in cloud service API specification languages and thus not enforced in implementations.*

### C. Implementation Defects

Lastly, we observe discrepancies caused by implementation defects, including unimplemented features and implementation bugs in the emulators and the cloud services.

1) *Unimplemented features in emulators:* A significant percentage of discrepancies are due to unimplemented features in emulators, accounting for 18% in Azure Storage, 16% in S3, and 74% in AWS DynamoDB. The emulators’ responses to these unimplemented APIs vary. For example, Azurite responds with a 500 error and the message “*Current API is not implemented yet,*” which leads to four wasted retries by the SDK. Whereas the AWS emulator issues a 400 error, as shown in Figure 5a, without triggering retries on the client side. According to the coverage reference of LocalStack [31], there are 16 (30%) unimplemented DynamoDB APIs and seven (7%) unimplemented S3 APIs. Azurite states that more features will be supported based on the needs of customers [54].

It is expensive to implement and maintain the large number of cloud service APIs (with high fidelity) in the emulator. Hence, existing emulators take a utility-driven approach to only support commonly used APIs (Figure 2). However, if a project relies on unimplemented APIs, the limited support becomes an obstacle for emulator-based testing.

2) *Emulator bugs:* Bugs were identified as the root causes of 15% of Azure Storage, 16% of S3, and 4% of DynamoDB discrepancies. Ten emulator bugs were found across the three classes of services (three in Storage, six in S3, and one in DynamoDB). For example, in Figure 5b, during tests involving AWS S3’s object restoration API, we encountered different responses to invalid keys. While the cloud service correctly rejected invalid keys with a 404 NoSuchKey error, the emulator returned a 500 Internal Error due to a bug that attempted to access a non-existent “`storage_class`” attribute.

3) *Cloud service bugs:* We also find two bugs in the cloud service that resulted in inconsistencies with the emulator. As shown in Figure 5c, specifying a lease duration for a blob

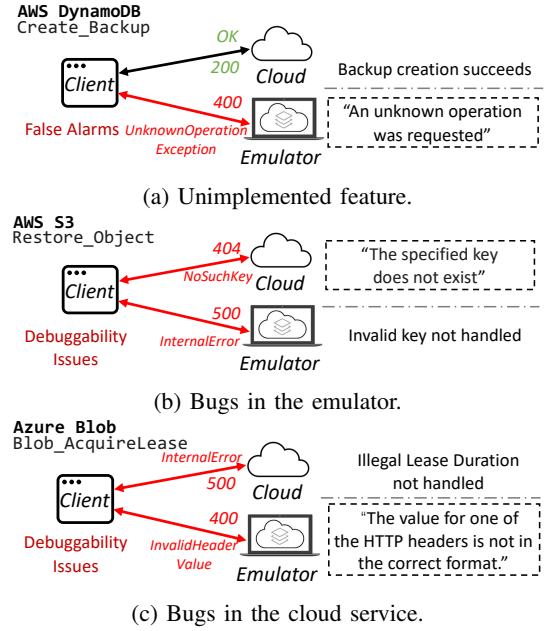


Fig. 5: Discrepancies caused by implementation defects.

outside the documented range of 15 to 60 seconds, particularly with an excessively large value, led to a 500 Internal Server Error in the cloud service. Contrastingly, the emulator appropriately responded with a 400 InvalidHeaderValue error, correctly identifying the lease duration as invalid. A similar bug is found in the API for acquiring a container lease.

**Finding 7.** *With the current development of cloud emulators as a reactive practice to cloud services, discrepancies due to unimplemented features and bugs would largely continue, even with active bug fixing and feature requests. Resolving discrepancies needs novel mitigation techniques.*

### D. Interactions with Developers

We contacted developers of two studied projects; the developers considered service-emulator discrepancies problematic. Currently, their practice is to always run tests that invoke discrepant APIs with actual cloud services, which is not ideal. Unfortunately, from software developers’ perspective, there is not much they could do as consumers of these services.

Therefore, we interacted with developers of the two emulators (Azurite and LocalStack). The reported discrepancies were appreciated. Six of them (caused by implementation bugs) have been confirmed (five fixed) as discussed in §V-C2. None reported discrepancy was rejected. We also reported bugs of cloud services; both were confirmed (see §V-C3).

## VI. DISCREPANCY MITIGATION: A DISCUSSION

It is easier to ask for more specifications (§V-A and §V-B) and faster bug fixes (§V-C). However, it is harder to fundamentally eliminate all the aforementioned discrepancies, as many of them are rooted in the essential complexity of software engineering as well as today’s common practices. We discuss



a few arguably radical ideas or new practices, hopefully to shed light on viable directions to addressing discrepancies.

#### A. An Active Role of Cloud Service Providers

Our fuzzing-based differential testing shows the effectiveness of detecting discrepancies between cloud services and emulators. Cloud service providers can adopt similar practices; they can run the differential testing continuously upon code changes of the emulators or cloud service implementations. Note that cloud providers already run REST API fuzzing to find security bugs in cloud backend implementations [4], [21], [38]. As service providers have more resources and insights into the implementations, they are in a better position than researchers or application developers to discover discrepancies and should take an active role in communicating and documenting them (e.g., as a part of the API specification). Certainly, maintenance of discrepancy-related documents can be challenging and the cost cannot be overlooked.

**Incentives.** Our discussion with the cloud providers show strong incentives to address discrepancies and improve emulator fidelity. High-fidelity service emulation would improve developer experience and help promote adoption of cloud services, which is why providers offer official emulators. Such benefits make a better strategy than forcing customers to run all their tests with the cloud services. Certainly, there is no free lunch and it may take major efforts to enforce conformance and achieve high fidelity; hence, the incentive structure can be complicated. We did not hear concerns regarding confidentiality when specifying external semantics and behavior of APIs.

#### B. Formal Models as Emulators

Essentially, discrepancies are introduced through the current practice of implementing emulators. Our private communication with a major cloud service provider tells us that the emulators are often not developed by the same engineering team that developed the cloud services, and the emulators are developed *reactive* to the cloud services; for third-party emulators like LocalStack, it is unavoidable. So, without comprehensive formal specifications, discrepancies are inevitable.

One way to resolve discrepancies is to change how emulators are built today. We envision the use of *executable* formal models of cloud services as the emulators. Essentially, the emulator, as the formal model, defines the specifications of the cloud service implementations, which are rigorously tested or verified for compliance. Recent efforts from Amazon [55] show the promise of developing executable reference models as specifications to be checked against the implementation of ShardStore, a key-value storage node of Amazon S3. Similar efforts have been made for other system domains [10], [56], [57]. In principle, these models can further be developed into first-class emulators for application testing.

#### C. “POSIX” for Cloud Service APIs

One fact that makes cloud emulators particularly prone to discrepancies is the lack of a standard such as POSIX for operating system call APIs. Today, cloud service providers

expose APIs with different semantics, constraints, and error codes, even for the same types of services. Without a standardized API, cloud and emulator developers must navigate diverse semantics and error handling for even similar services within the same or across different platforms. As a result, implementations of the APIs, whether by the cloud services or emulators, tend to be error-prone and inconsistent. As a result, we believe that a unified API standard would effectively reduce discrepancies in practice. With the incentives from sky computing [58] and hybrid cloud [59], [60], such a unified API standard may be possible.

#### D. Economic Cloud Services for Testing

With the prevalent discrepancies (§IV), testing of cloud-based software would have to largely rely on cloud services. To reduce cost, one can minimize the frequency of running tests with real cloud services (e.g., only do so before deployment, not for CI). If cost is the main concern (a recent survey [61] shows that cost is a major barrier to cloud service adoption), one solution is to provide cheap cloud services for testing. The high cost of cloud services is often driven by pursuits for performance using powerful hardware and fault tolerance using redundancy (§II-A). But, functional and correctness testing may not need either of them. We envision low-cost cloud services specified for software testing (not for production), with ideas such as using dated hardware [62] and cheap, renewable energy for intermittent services [63]. Certainly, low-cost services do not address other needs of emulators, such as convenience and hermetic environments [64].

#### E. Hybrid Cloud-Emulator Testing

One principle to mitigate discrepancies without exclusively using cloud services for testing is to acknowledge imperfect emulators and make the best use of them—selectively running tests on emulators when the emulation is not discrepant and on cloud services otherwise. We term such an approach *hybrid cloud-emulator testing* and explore it in §VII.

### VII. HYBRID CLOUD-EMULATOR TESTING

To evaluate the effectiveness of hybrid cloud-emulator testing as a short-term discrepancy mitigation (§VI-E), we developed a tool named ET that determines whether a test should be run with emulators or cloud services. The principle is to run discrepant tests with cloud services for safety while running the remaining tests with the emulators for efficiency. Note that it is hard to selectively use emulators and cloud services within a test without expensive state synchronization.

#### A. Policies

ET supports three different but complementary policies:

##### 1) Selection by discrepant APIs (API-based selection):

This policy assumes apriori knowledge of the set of discrepant APIs for a given service (which requires maintenance of the discrepant API set; see §VI-A). It first runs all tests on the emulator and monitors their REST API calls using a local proxy, as in [65]. If a test invokes a known discrepant API, its result is discarded, and the test is rerun on cloud service.

2) *Selection by in-situ API monitoring (Monitoring-based selection)*: API-based selection (§VII-A1) assumes having accurate, comprehensive discrepancy information apriori, which can be costly (§VI-A) and is often incomplete in practice (as shown in §IV-A, discrepancies are found in developer-certified APIs). However, without knowledge of discrepancies, all the tests have to run with cloud services.

ET supports a new monitoring-based policy that offloads certain tests from the cloud services to the emulators to reduce cost. The high-level idea is to maintain a “safe list” database of API call sequences. Each API call in a sequence consists of the API ID, the call’s request (with parameters), and the response. ET starts with an empty safe list. For each test, it first runs the test on the emulator and monitors the API call sequences. If the sequence is not present in the safe list, the test is assumed to be discrepant and is rerun on the cloud service. The API sequence is added to the safe list if its result from the emulator matches that from the cloud service.

On the other hand, if the API sequence is found in the safe list, then ET skips running the test with the cloud services, with the rationale that the fidelity of interactions has already been validated by a real cloud-based test run, and thus can be saved. Figure 6 illustrates the workflow.

Note that the analysis considers the entire API call sequence instead of individual APIs (Figure 4 shows an example where the discrepancy only manifests with specific sequence). ET serializes the API calls at its local proxy. For fast comparison, each sequence is stored as an ordered list of hashes while each hash represents an API along with its parameters and response; collisions are chained upon occurrence. We implement masks to exclude intrinsically non-deterministic parameters and fields in the response, such as timestamps.

The monitoring-based selection accounts for nondeterminism of test execution due to multi-threading and event-based asynchrony. The essence of the policy is to validate external behavior of API calls issued by the test on the emulator with real cloud services. ET does not make assumptions on the internal implementation of test code or system under test.

Despite no apriori discrepancy knowledge, the monitoring-based selection policy can outperform the API-based policy in §VII-A1 in certain cases, because it performs a fine-grained, parameter-level analysis, instead of labeling the entire API as in §VII-A1. As shown in §IV-C, discrepancies typically manifest via specific parameter values rather than universally across the API; a discrepant test can still be run on the emulator if it only uses the API with “safe” parameters and does not manifest discrepant results. On the other hand, as any unseen API call sequence is considered unsafe under this policy, the effectiveness of the monitoring-based selection relies on the coverage of the recorded safe sequences.

3) *Combined selection policy*: ET also supports a combined policy that integrates the API-based and monitoring-based selection policies to take the advantages of both policies. It has the same assumption of an apriori set of discrepant APIs, as in §VII-A1. Like the API-based selection (§VII-A1), a test first runs on the emulator; if the test is not a discrepant test,

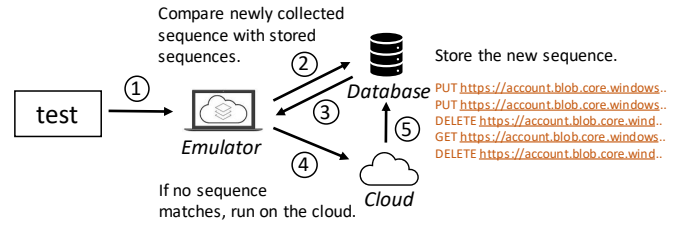


Fig. 6: Workflow of the monitoring-based selection (§VII-A2).

TABLE VII: Savings of cloud service API invocations with different policies of ET (averaged across five commits)

| Project            | Total Tests | Total Requests | # Saved Requests |         |         |
|--------------------|-------------|----------------|------------------|---------|---------|
|                    |             |                | §VII-A1          | §VII-A2 | §VII-A3 |
| Orleans            | 189         | 117,905        | 29.4%            | 0.4%    | 29.6%   |
| Insights           | 171         | 5,249          | 4.3%             | 47.2%   | 47.2%   |
| Durabletask        | 101         | 79,654         | 0%               | 0%      | 0%      |
| Streamstone        | 75          | 590            | 0%               | 99.2%   | 99.2%   |
| IdentityAzureTable | 51          | 9,860          | 100%             | 0.4%    | 100%    |

then it does not need to rerun on the cloud services. So, non-discrepant tests are always saved and are not affected by the coverage limitation of the monitoring-based selection. Only if a test is a discrepant test, we apply in-situ discrepant analysis as in §VII-A2 to conduct more fine-grained, parameter-level policy. Hence, the combined policy is always more effective than the API- and monitoring-based policies individually.

## B. Evaluation

We evaluate the three policies in §VII-A in terms of cost savings measured by the number of calls to cloud APIs. We assume a continuous integration (CI) setup as the monitoring-based selection benefits continuous testing the most (its benefit is correlated with the comprehensiveness of the safe list).

We select five projects that use Azure APIs (ET currently only supports .NET applications) and use all the related tests (Table II). We select the five projects with the most tests that invoke Azure APIs. We evaluate ET with the most recent five commits to simulate CI<sup>2</sup> and record the cost saving for each commit (we run only five commits due to the constraint of our cloud education credits). All the tests will be run for each commit that changes system code or test code. Note that regression test selection [66]–[68] does not apply to those tests which are not unit tests but mostly integration and system tests.

In the evaluation, we use the same versions of the Azurite emulator and Azure services [4], without considering their software evolution (which may lead to test regressions).

**Results.** As shown in Table VII, ET effectively reduces the amount of invocations to cloud APIs. Interestingly, the three policies bring different benefits across projects, with the combined policy (§VII-A3) achieving the most cost savings.

The API-based selection (§VII-A1) achieves substantial savings for two out of five projects. Specifically, it achieves

<sup>2</sup>We assume the maintenance of an always updated set of discrepant APIs apriori for the API-based and the combined policies (§VII-A1 and §VII-A3).

a 100% saving for IdentityAzureTable where none of its tests issues discrepant APIs. However, it achieves no saving for Streamstone and DurableTask, because all their tests issue at least one discrepant API.

The monitoring-based selection (§VII-A2) achieves substantial savings for two different projects (Insights and Streamstone) but not the others. Our investigation reveals that the effectiveness of this policy largely depends on the ordering determinism of API call sequences. Since tests that invoke cloud APIs are typically large system/integration tests with large numbers of API calls, the sequences recorded during the tests on five commits are insufficient. We expect that a longer continuous testing process may increase the benefit.

**Finding 8.** *ET shows that by selectively running tests on emulators, it is promising to reduce the cost of cloud-based software testing, in terms of the cost of calling cloud service APIs, while achieving high-fidelity testing.*

## VIII. THREATS TO VALIDITY

Our study is based on the five cloud services (three Azure services and two AWS services) and two emulators (Azurite and LocalStack). We believe that the studied cloud services and emulators are representative, but our results may not generalize to other cloud services, especially those using different practices of API design and specification. For example, the issues could be more severe for smaller providers. It would be interesting to check the gaps among different providers from different tiers, which is our future work. We recommend readers to focus on overall trends and not on precise number.

Similarly, our analysis of discrepancy impacts (§IV-C) and ET’s evaluation results (§VII-B) are based on existing test suites of a few cloud-based projects. They may not generalize to other projects as they depend on API usage characteristics of projects and their tests (e.g., invoked APIs and frequencies). In principle, projects that use cloud services more extensively face higher impacts of discrepancies. With wider adoption of cloud-based programming, we expect cloud-emulator discrepancies to be common issues for software testing.

The discrepancies analyzed in this paper are limited to the black-box SDK API fuzzer we developed based on RESTler (§III), and we do not claim completeness of studied discrepancies. A more powerful fuzzer, especially a white-box one, may cover more discrepant APIs. As a best effort, we run our fuzzer for more than ten hours against each studied emulator and stop the fuzzing when we do not observe any new discrepancies.

Lastly, we are not concerned with faults that occur during the API invocations, such as timeout due to network delays. Recent work [65] shows that timeouts on both the request and response paths of a REST API invocation can reveal different behaviors, which we would like to study as future work.

## IX. RELATED WORK

**REST API fuzzing.** Recent work has developed advanced REST API fuzzing techniques to test web and cloud services, with the goal of finding bugs and vulnerabilities in web service

implementations [4], [21], [38], [69]–[78]. Differently, this paper focuses on the software projects that use cloud services, instead of the backend implementations of cloud services. Our goal is to understand discrepancies between the emulator and the cloud services and their implications for software testing.

We developed our fuzzer based on the fuzzing approach of RESTler [21]. As discussed in §III, we did not directly use RESTler (or other REST API fuzzers) because most projects only interact with SDK APIs, not REST APIs, and REST API fuzzers generate API calls that would not be output by SDKs.

**Fidelity of emulation.** Prior work has studied the fidelity of emulation environments in other domains, such as honeypots for security analysis [79]–[81]. The closest related work is the research on the fidelity of emulated execution environments such as virtual devices for mobile app testing [82]–[84] and the efforts of building usable, effective mobile emulators [85], [86]. The goal is to maximize app testing on emulated devices and minimize testing on real physical mobile devices (which are more expensive and hard to manage [82]).

Our work shares similar high-level goals and tradeoffs (cost-efficiency versus safety). But, we address a different fidelity problem raised by the emerging cloud-based programming model. The discrepancies are not due to deficiency or incompleteness of device emulation but are rooted in inconsistent implementations of weakly specified APIs.

**Backward compatibility.** The studied discrepancies are different from backward incompatibility studied in prior work [4], [87]–[90]. We do not study the evolution of emulators or cloud service APIs in this paper, though certain discrepancies can be caused by regression [4]. There are also studies on mock libraries for unit tests [91]–[93]; few of them concern fidelity of mock objects—unlike emulation, mocking is not expected to provide fidelity but offers a way to control external APIs.

## X. CONCLUSION

With the rise of cloud-based programming models, testing cloud-based software safely and cost-efficiently becomes a challenge. This paper analyzes the fidelity of commonly used cloud emulators for software development and offline testing. Our results show that discrepancies between the emulator and the cloud services are prevalent today, affecting the safety and trustworthiness of testing results of cloud-based software. We discuss both accidental and essential root causes of these discrepancies and envision new practices and techniques to mitigate them, even though fundamentally eliminating them can be hard. We also show the promise of leveraging imperfect emulators for cost-efficient testing.

## ACKNOWLEDGMENTS

We thank Darko Marinov, Hao Lin, Talha Waheed, and Owolabi Legunsen for their valuable discussions and feedback. We thank the LocalStack developers for helping us understand the implementation. We thank Marina Polishchuk for her help on RESTler. This work was supported in part by NSF CNS-2130560, CNS-2145295, and an IIDAI grant.

## REFERENCES

- [1] “Azure products,” <https://azure.microsoft.com/en-us/products>, 2024.
- [2] “Google Cloud products,” <https://cloud.google.com/products>, 2024.
- [3] “AWS Cloud Products,” <https://aws.amazon.com/products>, 2024.
- [4] P. Godefroid, D. Lehmman, and M. Polishchuk, “Differential Regression Testing for REST APIs,” in *ISSTA*, 2020.
- [5] “Azure Blob Storage Cost,” <https://azure.microsoft.com/en-us/pricing/details/storage/blobs/>, 2023.
- [6] “AWS Pricing Calculator,” <https://calculator.aws/#/?nc2=pr&refid=f42fef03-b1e6-4841-b001-c44b4eccaf41>, 2024.
- [7] John Micco, “The State of Continuous Integration Testing @Google,” in *ICST*, 2017.
- [8] C. Leong, A. Singh, M. Papadakis, Y. L. Traon, and J. Micco, “Assessing Transition-Based Test Selection Algorithms at Google,” in *ICSE-SEIP*, 2019.
- [9] S. Wang, X. Lian, D. Marinov, and T. Xu, “Test Selection for Unified Regression Testing,” in *ICSE*, 2023.
- [10] T. Ridge, D. Sheets, T. Tuerk, A. Giugliano, A. Madhavapeddy, and P. Sewell, “SibylFS: Formal Specification and Oracle-Based Testing for POSIX and Real-World File Systems,” in *SOSP*, 2015.
- [11] Rick Timmis, “Xiotech accelerates their development workflows on cloud using LocalStack!” <https://localstack.cloud/blog/2023-07-05-case-study-xiotech/>.
- [12] Microsoft Docs, “Use the Azurite emulator for local Azure Storage development,” <https://docs.microsoft.com/en-us/azure/storage/common/storage-use-azurite>, 2024.
- [13] “Introduction to Azure Storage,” <https://learn.microsoft.com/en-us/azure/storage/common/storage-introduction>, 2024.
- [14] azurite-1465, “Table storage having wrong constraint,” <https://github.com/Azure/Azurite/issues/1465>.
- [15] azurite-946, “Missing parameter in API response,” <https://github.com/Azure/Azurite/issues/946>.
- [16] azurite-5, “Requesting for a feature,” <https://github.com/Azure/Azurite/issues/5>.
- [17] nuget-insights-30, “List of issues which blocked CI tests for Insights,” <https://github.com/NuGet/Insights/issues/30>.
- [18] Localstack, “LocalStack – A fully functional local cloud stack,” <https://localstack.cloud/>, 2024.
- [19] Azure/durablerequest, <https://github.com/Azure/durablerequest>.
- [20] “Azurite Support Matrix,” <https://github.com/Azure/Azurite?tab=readme-ov-file#support-matrix>.
- [21] V. Atlidakis, P. Godefroid, and M. Polishchuk, “RESTler: Stateful REST API Fuzzing,” in *ICSE*, 2019.
- [22] “AWS S3 REST APIs,” [https://docs.aws.amazon.com/AmazonS3/latest/API/API\\_Operations.html](https://docs.aws.amazon.com/AmazonS3/latest/API/API_Operations.html), 2023.
- [23] “Azurite Swagger of Azure Blob Storage,” <https://github.com/Azure/Azurite/blob/main/swagger/blob-storage-2021-10-04.json>.
- [24] “OpenAPI Specification,” <https://www.openapis.org/>, 2024.
- [25] “Azure Blob Storage Cost Breakdown,” <https://azure.github.io/Storage/docs/application-and-user-data/code-samples/estimate-block-blob/>, 2023.
- [26] “Azure Locally-redundant storage (LRS),” <https://learn.microsoft.com/en-us/azure/storage/common/storage-redundancy#locally-redundant-storage>, 2024.
- [27] “Azure Zone-redundant storage (ZRS),” <https://learn.microsoft.com/en-us/azure/storage/common/storage-redundancy#zone-redundant-storage>, 2024.
- [28] “Azure Locally-redundant storage (GRS),” <https://learn.microsoft.com/en-us/azure/storage/common/storage-redundancy#geo-redundant-storage>, 2024.
- [29] Brian Harry, “The largest Git repo on the planet,” <https://devblogs.microsoft.com/bharry/the-largest-git-repo-on-the-planet/>, 2017.
- [30] R. Potvin and J. Levenberg, “Why Google Stores Billions of Lines of Code in a Single Repository,” *Communications of the ACM (CACM)*, vol. 59, no. 7, pp. 78–87, Jun. 2016.
- [31] “LocalStack Coverage,” <https://docs.localstack.cloud/references/coverage/>, 2023.
- [32] Azure/autorest, <https://github.com/Azure/autorest>, 2024.
- [33] “Regeneration Protocol Layer from Swagger by Autorest,” <https://github.com/Azure/Azurite/blob/main/CONTRIBUTION.md#regeneration-protocol-layer-from-swagger-by-autorest>, 2024.
- [34] “LocalStack Weekly ASF Update Workflow,” <https://github.com/localstack/localstack/blob/master/.github/workflows/asf-updates.yml>, 2024.
- [35] “Mocking in Unit Tests,” <https://microsoft.github.io/code-with-engineering-playbook/automated-testing/unit-testing/mocking>, 2023.
- [36] M. Fowler, “Mocks Aren’t Stubs,” <https://martinfowler.com/articles/mocksArentStubs.html>, 2007.
- [37] Statista, “Amazon Maintains Cloud Lead as Microsoft Edges Closer,” <https://www.statista.com/chart/18819/worldwide-market-share-of-leading-cloud-infrastructure-service-providers/>, 2024.
- [38] V. Atlidakis, P. Godefroid, and M. Polishchuk, “Checking Security Properties of Cloud Service REST APIs,” in *ICST*, 2020.
- [39] LocalStack, “Coverage S3 - LocalStack Documentation,” [https://docs.localstack.cloud/references/coverage/coverage\\_s3/](https://docs.localstack.cloud/references/coverage/coverage_s3/), 2023.
- [40] LocalStack, “Coverage DynamoDB - LocalStack Documentation,” [https://docs.localstack.cloud/references/coverage/coverage\\_dynamodb/](https://docs.localstack.cloud/references/coverage/coverage_dynamodb/), 2023.
- [41] LocalStack, “Parity Testing - LocalStack Documentation,” <https://docs.localstack.cloud/contributing/parity-testing/>, 2023.
- [42] “Azure Storage API: ContainerRestore,” <https://github.com/Azure/Azurite/blob/d544d16f910e490fd9db5565459df701895308f/swagger/blob-storage-2021-10-04.json#L1554>, 2024.
- [43] “Azure Storage API: BlockBlob-StageBlockFromURL,” <https://github.com/Azure/Azurite/blob/d544d16f910e490fd9db5565459df701895308f/swagger/blob-storage-2021-10-04.json#L7069>, 2024.
- [44] Q. Luo, F. Hariri, L. Eloussi, and D. Marinov, “An Empirical Analysis of Flaky Tests,” in *FSE*, 2014.
- [45] O. Parry, G. M. Kapfhammer, M. Hilton, and P. McMinn, “A Survey of Flaky Tests,” *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 31, no. 17, pp. 1–74, Oct. 2021.
- [46] “Lease Blob - Microsoft Documentation,” <https://learn.microsoft.com/en-us/rest/api/storageservices/lease-blob?tabs=microsoft-entra-id>.
- [47] K. Lazar, M. Vetzler, G. Uziel, D. Boaz, E. Goldbraich, D. Amid, and A. Anaby-Tavor, “SpecCrawler: Generating OpenAPI Specifications from API Documentation Using Large Language Models,” *arXiv:2402.11625*, 2024.
- [48] M. Endres, S. Fakhoury, S. Chakraborty, and S. K. Lahiri, “Can Large Language Models Transform Natural Language Intent into Formal Method Postconditions?” in *FSE*, 2024.
- [49] A. Decrop, G. Perrouin, M. Papadakis, X. Devroey, and P.-Y. Schobbens, “You Can REST Now: Automated Specification Inference and Black-Box Testing of RESTful APIs with Large Language Models,” *arXiv:2402.05102*, 2024.
- [50] L. Ma, S. Liu, Y. Li, X. Xie, and L. Bu, “SpecGen: Automated Generation of Formal Program Specifications via Large Language Models,” *arXiv:2401.08807*, 2024.
- [51] R. Yandrapally, S. Sinha, R. Tzoref-Brill, and A. Mesbah, “Carving UI Tests to Generate API Tests and API Specification,” in *ICSE*, 2023.
- [52] R. Huang, M. Motwani, I. Martinez, and A. Orso, “Generating REST API Specifications through Static Analysis,” in *ICSE*, 2024.
- [53] C. Hathhorn, C. Ellison, and G. Roşu, “Defining the Undefinedness of C,” in *PLDI*, 2015.
- [54] “Azurite README,” <https://github.com/Azure/Azurite/blob/main/README.md>.
- [55] J. Bornholt, A. Kaufmann, J. Li, A. Krishnamurthy, E. Torlak, and X. Wang, “Specifying and Checking File System Crash-Consistency Models,” in *ASPLOS*, 2016.
- [56] X. Sun, W. Ma, J. T. Gu, Z. Ma, T. Chajed, J. Howell, A. Lattuada, O. Padon, L. Suresh, A. Szekeres, and T. Xu, “Anvil: Verifying Liveness of Cluster Management Controllers,” in *OSDI*, 2024.
- [57] S. Bishop, M. Fairbairn, M. Norrish, P. Sewell, M. Smith, and K. Wansbrough, “Rigorous Specification and Conformance Testing Techniques for Network Protocols, as applied to TCP, UDP, and Sockets,” in *SIGCOMM*, 2005.
- [58] S. Chasins, A. Cheung, N. Crooks, A. Ghodsi, K. Goldberg, J. E. Gonzalez, J. M. Hellerstein, M. I. Jordan, A. D. Joseph, M. W. Mahoney, A. Parameswaran, D. Patterson, R. A. Popa, K. Sen, S. Shenker, D. Song, and I. Stoica, “The Sky Above the Clouds: A Berkeley View on the Future of Cloud Computing,” *arXiv:2205.07147*, 2022.
- [59] Google Cloud, “What is a Hybrid Cloud?” <https://cloud.google.com/learn/what-is-hybrid-cloud>.
- [60] IBM Hybrid Cloud, “Hybrid cloud solutions,” <https://www.ibm.com/hybrid-cloud>.

- [61] M. Loukides, "The Cloud in 2021: Adoption Continues," O'Reilly Media, Tech. Rep., 2021.
- [62] J. Wang, U. Gupta, and A. Sriraman, "Giving Old Servers New Life at Hyperscale," in *HotInfra*, 2023.
- [63] P. Ambati, I. Goiri, F. Frujeri, A. Gun, K. Wang, B. Dolan, B. Corell, S. Pasupuleti, T. Moscibroda, S. Elnikety, M. Fontoura, and R. Bianchini, "Providing SLOs for Resource-Harvesting VMs in Cloud Platforms," in *OSDI*, 2020.
- [64] C. Narla and D. Salas, "Hermetic Servers," <https://testing.googleblog.com/2012/10/hermetic-servers.html>, Oct. 2012.
- [65] Y. Chen, X. Sun, S. Nath, Z. Yang, and T. Xu, "Push-Button Reliability Testing for Cloud-Backed Applications with Rainmaker," in *NSDI*, 2023.
- [66] G. Rothermel and M. Harrold, "Analyzing Regression Test Selection Techniques," *IEEE Transactions on Software Engineering (TSE)*, vol. 22, no. 8, pp. 529–551, Aug. 1996.
- [67] G. Rothermel and M. Harrold, "A Safe, Efficient Regression Test Selection Technique," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 6, no. 2, p. 173–210, Apr. 1997.
- [68] S. Yoo and M. Harman, "Regression Testing Minimisation, Selection and Prioritisation: A Survey," *Software Testing, Verification & Reliability*, vol. 22, no. 2, p. 67–120, Mar. 2012.
- [69] C. Lyu, J. Xu, S. Ji, X. Zhang, Q. Wang, B. Zhao, G. Pan, W. Cao, P. Chen, and R. Beyah, "MINER: A Hybrid Data-Driven Approach for REST API Fuzzing," in *USENIX Security*, 2023.
- [70] P. Godefroid, B.-Y. Huang, and M. Polishchuk, "Intelligent REST API Data Fuzzing," in *ESEC/FSE*, 2020.
- [71] O. Bastani, R. Sharma, A. Aiken, and P. Liang, "Synthesizing Program Input Grammars," in *PLDI*, 2017.
- [72] P. Godefroid, H. Peleg, and R. Singh, "Learn&Fuzz: Machine Learning for Input Fuzzing," in *ASE*, 2017.
- [73] H. Wu, L. Xu, X. Niu, and C. Nie, "Combinatorial Testing of RESTful APIs," in *ICSE*, 2022.
- [74] J. C. Alonso, "Automated Generation of Realistic Test Inputs for Web APIs," in *ESEC/FSE*, 2021.
- [75] A. Martin-Lopez, S. Segura, and A. Ruiz-Cortés, "RESTTest: Automated Black-Box Testing of RESTful Web APIs," in *ISSTA*, 2021.
- [76] A. Arcuri, "RESTful API Automated Test Case Generation with EvoMaster," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 28, no. 1, pp. 1–37, Jan. 2019.
- [77] E. Viglianisi, M. Dallago, and M. Ceccato, "RESTTESTGEN: Automated Black-Box Testing of RESTful APIs," in *ICST*, 2020.
- [78] Y. Liu, Y. Li, G. Deng, Y. Liu, R. Wan, R. Wu, D. Ji, S. Xu, and M. Bao, "Morest: Model-based RESTful API Testing with Execution Feedback," in *ICSE*, 2022.
- [79] F. Dang, Z. Li, Y. Liu, E. Zhai, Q. A. Chen, T. Xu, Y. Chen, and J. Yang, "Understanding Fileless Attacks on Linux-based IoT Devices with HoneyCloud," in *MobiSys*, 2019.
- [80] C. Kreibich and J. Crowcroft, "Honeycomb – Creating Intrusion Detection Signatures Using Honey Pots," in *SIGCOMM-CCR*, 2004.
- [81] M. D. Vrabie, J. T. Ma, J. Chen, D. Moore, E. Vandekieft, A. C. Snoeren, G. M. Voelker, and S. Savage, "Scalability, Fidelity, and Containment in the Potemkin Virtual Honeyfarm," in *SOSP*, 2005.
- [82] H. Lin, J. Qiu, H. Wang, Z. Li, L. Gong, D. Gao, Y. Liu, F. Qian, Z. Zhang, P. Yang, and T. Xu, "Virtual Device Farms for Mobile App Testing at Scale: A Pursuit for Fidelity, Efficiency, and Accessibility," in *MobiCom*, 2023.
- [83] H. Lin, J. Qiu, H. Wang, Z. Li, L. Gong, D. Gao, Y. Liu, F. Qian, Z. Zhang, P. Yang, and T. Xu, "Take the Blue Pill: Pursuing Mobile App Testing Fidelity, Efficiency, and Accessibility with Virtual Device Farms," *SIGMOBILE Mobile Computing and Communications (GetMobile)*, vol. 28, no. 1, pp. 5–9, Mar. 2024.
- [84] H. Cai, Z. Zhang, L. Li, and X. Fu, "A Large-Scale Study of Application Incompatibilities in Android," in *ISSTA*, 2019.
- [85] J. Qiu, Z. Zhou, Y. Li, Z. Li, F. Qian, H. Lin, D. Gao, H. Su, X. Miao, Y. Liu, and T. Xu, "vSoC: Efficient Virtual System-on-Chip on Heterogeneous Hardware," in *SOSP*, 2024.
- [86] D. Gao, H. Lin, Z. Li, C. Huang, L. Gong, F. Qian, Y. Liu, and T. Xu, "Trinity: High-Performance Mobile Emulation through Graphics Projection," in *OSDI*, 2022.
- [87] L. Chen, F. Hassan, X. Wang, and L. Zhang, "Taming Behavioral Backward Incompatibilities via Cross-Project Testing and Analysis," in *ICSE*, 2020.
- [88] Y. Zhang, J. Yang, Z. Jin, U. Sethi, K. Rodrigues, S. Lu, and D. Yuan, "Understanding and Detecting Software Upgrade Failures in Distributed Systems," in *SOSP*, 2021.
- [89] C. Zhu, M. Zhang, X. Wu, X. Xu, and Y. Li, "Client-Specific Upgrade Compatibility Checking via Knowledge-Guided Discovery," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 32, no. 4, pp. 1–31, May 2023.
- [90] Y. Zhao, L. Li, K. Liu, and J. C. Grundy, "Towards Automatically Repairing Compatibility Issues in Published Android Apps," in *ICSE*, 2022.
- [91] D. Saff, S. Artzi, J. H. Perkins, and M. D. Ernst, "Automatic Test Factoring for Java," in *ASE*, 2005.
- [92] X. Wang, L. Xiao, T. Yu, A. Woepse, and S. Wong, "An Automatic Refactoring Framework For Replacing Test-Production Inheritance by Mocking Mechanism," in *ESEC/FSE*, 2021.
- [93] A. Arcuri, G. Fraser, and J. P. Galeotti, "Automated Unit Test Generation for Classes with Environment Dependencies," in *ASE*, 2014.