# GraphQLer: Enhancing GraphQL Security with Context-Aware API Testing

Omar Tsai*
omar@ztasecurity.com
Simon Fraser University
Burnaby, BC, Canada

Jianing Li
jianing_li_2@sfu.ca
Simon Fraser University
Burnaby, BC, Canada

Tsz Tung Cheung
ttc9@sfu.ca
Simon Fraser University
Burnaby, BC, Canada

Lejing Huang
lha102@sfu.ca
Simon Fraser University
Burnaby, BC, Canada

Hao Zhu
hza164@sfu.ca
Simon Fraser University
Burnaby, BC, Canada

Jianrui Xiao
jxa68@sfu.ca
Simon Fraser University
Burnaby, BC, Canada

Iman Sharafaldin
i.sharafaldin@fwdsec.com
Forward Security
Vancouver, BC, Canada

Mohammad A. Tayebi
tayebi@sfu.ca
Simon Fraser University
Burnaby, BC, Canada

## Abstract

GraphQL is an open-source data query and manipulation language designed for web applications, offering a flexible alternative to RESTful APIs. However, its dynamic execution model and lack of built-in security mechanisms introduce significant vulnerabilities, including unauthorized data access, denial-of-service (DoS) attacks, and injection threats. Existing GraphQL testing frameworks primarily focus on functional correctness, neglecting security risks arising from query interdependencies and execution context, leading to inadequate vulnerability detection. This paper introduces GraphQLer—the first context-aware security testing framework for GraphQL APIs. GraphQLer analyzes intricate relationships among mutations, queries, and objects by constructing a dependency graph, capturing critical security-relevant interdependencies. It intelligently chains associated queries and mutations to uncover authentication and authorization flaws, access control bypasses, and resource misuse. Furthermore, GraphQLer can track the use of internal resources in requests to detect data leakage, privilege escalation, and replay attack vectors. We evaluate GraphQLer's effectiveness through various testing scenarios on different GraphQL APIs. In terms of testing coverage, GraphQLer achieves a significant improvement, with an average coverage increase of 35%, and in some instances, a remarkable 84% boost compared to the best-performing baseline method. This enhanced coverage is achieved in less time than the baseline methods, making GraphQLer particularly valuable for time-constrained use cases. Additionally, GraphQLer demonstrates its capability in detecting a previously reported CVE and potential vulnerabilities within large-scale production APIs. These findings highlight GraphQLer's potential to proactively strengthen GraphQL API security, offering a robust and automated approach to identifying security weaknesses.

## 1 Introduction

In recent years, there has been a noticeable increase in the use of modular service-based methodologies in software development [24]. This paradigm shift has prompted many developers to rely on Application Programming Interfaces (APIs) to facilitate data exchange

---

*Corresponding author.

between different architectural modules [45]. While established API architectures like Representational State Transfer (REST) [29] and Simple Object Access Protocol (SOAP) [23] are well-known, the Graph Query Language (GraphQL), due to its relatively recent emergence, remains one of the least explored alternatives despite its usefulness and growing popularity [50].

GraphQL is an API specification that allows users to retrieve data from multiple sources efficiently using a single query. As of 2024, over 60% of developers are using GraphQL to expose APIs for their applications and personal websites [7]. GraphQL's core approach centers around structuring data and requests in a graph, enabling developers to explicitly define dependencies and simplify the retrieval of correlated data [26], [63]. Despite its advantages, GraphQL introduces unique security challenges that are often overlooked. Unlike traditional REST APIs, where endpoints and response structures are predefined, GraphQL allows clients to dynamically construct queries, increasing the risk of over-fetching, information disclosure, and unauthorized access [40].

In our extensive domain analysis, we have identified a significant gap in the availability of effective solutions for GraphQL security testing. Many existing GraphQL frameworks test APIs by exhaustively enumerating available methods, applying parameter fuzzing techniques, and subsequently sending requests to detect errors in the response. Current GraphQL testing tools, such as Zed Attack Proxy (ZAP) [47] and BurpSuite [31], primarily focus on detecting simple vulnerabilities like injection flaws or schema misconfigurations. However, these tools fail to capture complex security risks that arise from endpoint dependencies.

Other proposed methods utilize static analysis, ignoring the dynamic responses of the API [40], [22], [39], [61], [65], [66]. While these techniques detect basic bugs, they overlook chained vulnerabilities that arise from implicit dependencies between queries, mutations, and objects. Some improvements have been made using reinforcement learning to dynamically learn from the responses of the API; none have yet looked at the dependent nature of the endpoint dependencies [53]. Existing works focus on input variation but fail to model how query interactions can be exploited for privilege escalation, authorization bypass, or data extraction. As a result, GraphQL APIs remain vulnerable to sophisticated multi-step

attacks that existing solutions cannot detect. In summary, both established testing tools and academic studies predominantly focus on input and output relationships in GraphQL, largely overlooking the untapped potential of utilizing GraphQL dependencies to sequence test cases.

To address these critical shortcomings, we introduce GraphQLer, the first *automated context–aware* security testing solution for GraphQL APIs. The term "context–aware" emphasizes GraphQLer's unique ability to comprehend the complex relationships between input and output within GraphQL, while the term "automated" highlights its ability to generate and execute security-focused test requests without manual intervention. Our approach includes the following key components: 1) *Dependency inference*: We infer dependencies among objects, queries, and mutations within GraphQL to identify potential attack vectors, such as unauthorized access via transitive relationships; 2) *Search algorithms*: We employ effective search techniques to dynamically explore vulnerability chains that existing testing approaches overlook; and 3) *Resource reusability*: Our methodology allows for the dynamic reuse of re- sources generated during GraphQL API testing. Unlike existing GraphQL testing frameworks that rely on arbitrary input modifications, GraphQLer leverages real-time API interactions to construct attack scenarios based on actual API behavior.

To evaluate the effectiveness of GraphQLer, we perform extensive testing on various GraphQL APIs, including public APIs [64] and large-scale platform APIs. Our results demonstrate that GraphQLer significantly improves test coverage, with an average increase of 35% and, in some cases, a notable rise of 84% compared to the most effective baseline method. Furthermore, GraphQLer successfully identifies multiple security vulnerabilities, including a previously reported CVE and several new potential exploits that were undetected by existing tools. These findings highlight the importance of incorporating context-aware dependency analysis into GraphQL testing to effectively detect complex vulnerabilities. In summary, we present the following key contributions.

◇ We advance the state-of-the-art in GraphQL security testing by introducing a novel context-aware approach that leverages query and mutation dependencies to uncover security vulnerabilities.

◇ We present GraphQLer, a fully functional, open-source tool [1]. As an open-source project, GraphQLer enables community-driven improvements, making it applicable across diverse GraphQL implementations.

◇ We perform extensive evaluation of GraphQLer on real-world GraphQL APIs, demonstrating its superiority over baseline methods in terms of both test coverage and security vulnerability detection. Notably, GraphQLer successfully identifies a documented CVE and several new security flaws, proving its effectiveness in uncovering critical GraphQL security risks.

This paper is structured as follows: In Section 2, we cover key concepts and related work. Section 3 outlines the proposed method, and Section 4 describes its implementation. Section 5 presents experimental results. Section 6 discusses implications and future directions. Finally, Section 7 concludes the paper.

---

[1]The source code for GraphQLer is hosted on Github at https://github.com/omar2535/GraphQLer.

## 2 Background & Related Work

In this section, we provide background on GraphQL APIs, discuss vulnerabilities inherent to them, and review prior work that addresses the problem we aim to solve.

### 2.1 GraphQL Fundamentals

GraphQL originated as an internal specification at Facebook in 2012 [6] and later transitioned into a public open-source project in 2015. Since then, it has been widely adopted by many well-known applications, including Instagram, Twitter, and Shopify.

GraphQL has the following features [51], [35], [56]: 1) *Data as a graph*: This enables all data to be fetched in a single request and prevents the under-fetching problem; 2) *Strongly typed*: GraphQL's schema defines data types (objects) and the possible methods for querying and mutating data over the API. This ensures predictable results and simplifies error handling. Clients must explicitly specify the fields of the data type they want in the output, minimizing over-fetching; and 3) *Single endpoint*: All requests are processed through a single, consistent endpoint (typically `/graphql` on the API server). To understand our approach to testing GraphQL APIs in Section 3, it is important to first grasp a few fundamental concepts of GraphQL, as discussed below.

In a GraphQL schema, two primary types of data are defined: *scalars* and *objects*. Scalars are primitive data types that represent a single atomic value. Standard GraphQL scalar types include `Int`, `Float`, `String`, `Boolean` and `ID`. Objects are user-defined entities that represent the structural blueprint and relationships between data. An object type encapsulates multiple fields, which can be scalars, objects, or lists of any type. This composition forms an interconnected structure between object types, creating the *graph* inherent to the API's purpose and enabling clients to traverse and retrieve interconnected data seamlessly.

Within the primary data types of GraphQL, there is the concept of nullability, which determines whether a field can contain a NULL value. This is often indicated by a `NON-NULL` flag. When a type is marked as `NON-NULL`, it means that the field in that type must always have a value (for example, the *name* scalar in a `User` object). To use scalars and objects, a user must interact with a GraphQL API through *queries* and *mutations*:

◇ *Queries* are operations that retrieve data from the server, similar to the GET request in traditional REST APIs. A specialized form of query, called *introspection query*, allows clients to retrieve detailed information about the GraphQL schema directly from the API.

◇ *Mutations* are operations that modify server-side data, akin to create, update or delete actions in a database.

Both queries and mutations enable clients to specify the fields of the returned type, ensuring that users receive exactly what they need. This minimizes the chances of over-fetching or under-fetching data, optimizing the efficiency of data retrieval.

### 2.2 GraphQL Vulnerabilities

While GraphQL offers powerful query capabilities beyond those of traditional REST APIs, it also introduces unique vulnerabilities that can pose significant security risks if not properly managed. These vulnerabilities, unique to GraphQL's design and functionality,

extend beyond typical web application issues—such as broken access controls, injection flaws, and misconfigurations [48], [46]—and create additional, GraphQL-specific attack surfaces. We categorize these vulnerabilities into three main categories:

◇ *Query Abuse Vulnerabilities:* GraphQL's flexibility, while advantageous, can be exploited to abuse query functionality, either to overwhelm the server or gain unauthorized insights into the schema. A notable example is the abuse of the GraphQL *introspection query*, which provides a detailed schema map, revealing all available types, fields and relationships. This schema exposure enables attackers to construct more precise and potentially harmful queries. GraphQL is also vulnerable to *Denial of Service (DoS)* attacks through deeply nested or excessively repetitive queries, which can exhaust server resources, leading to degraded performance or downtime.

◇ *Injection Vulnerabilities:* GraphQL applications are susceptible to various injection threats, which are among the most critical security risks. *SQL Injection* is a primary concern, where unsanitized inputs can be concatenated into SQL queries, allowing unauthorized data access or manipulation. Additionally, *Path Injection* and *Cross-Site Scripting (XSS)* attacks can occur when user inputs are improperly sanitized before being rendered on the client side, potentially leading to data theft, session hijacking or the distribution of malicious scripts. Poor handling of GraphQL inputs thus becomes a vector for injection attacks, compromising both the backend and end-users.

◇ *Access Control Vulnerabilities:* These vulnerabilities arise when GraphQL APIs fail to enforce adequate authorization policies, allowing unauthorized data access. Issues include *Insecure Direct Object References (IDOR)*, where attackers manipulate object identifiers to retrieve restricted data, and *batched attacks*, which combine multiple operations into a single request to bypass individual security checks. Addressing access control flaws is complex, requiring an understanding of dependencies in the GraphQL schema and the relationships among various fields, queries and mutations. Testing for these issues necessitates a context-aware, dependency-based testing approach that can accurately map schema relationships to simulate realistic user scenarios and uncover access control flaws that traditional testing methods often overlook.

## 2.3 Related works

The concept of dependency-aware testing is not new, as demonstrated by the work of various frameworks and testers proposed in the field of software validation [25], [52], [57], [43]. In the realm of REST API testing, dependency-aware methods like Microsoft's RESTler [21], foREST [42] and RestTestGen [62] are capable of comprehending the relationships between requests. However, dependency based testing for GraphQL remains very limited. To our knowledge, there is no existing literature on dependency-aware testing for GraphQL. One semantically-aware method proposed by Dodds et al. [36] addresses both REST and GraphQL APIs, yet its dependency-aware feature is solely built for fuzzing REST APIs, while the GraphQL fuzzer operates as a simple find-replace tester.

Since GraphQL functions as a specification layer over HTTP, several commercially available tools, such as BurpSuite [31] and

```
- name: createCurrency
  description: null
  args:
    - name: abbreviation
      description: null
      type:
        kind: NON_NULL
        name: null
        ofType:
          kind: SCALAR
          name: String
          ofType: null
      defaultValue: null
    - name: symbol
      description: null
      type:
        kind: NON_NULL
        name: null
        ofType:
          kind: SCALAR
          name: String
          ofType: null
      defaultValue: null
    - name: country
      description: null
      type:
        kind: SCALAR
        name: String
        ofType: null
      defaultValue: null
  type:
    kind: OBJECT
    name: Currency
    ofType: null
  isDeprecated: false
  deprecationReason: null
```

**Figure 1: Sample specification of a GraphQL mutation.**

ZAP [47], enable testing of a GraphQL server's generic responses. However, these tools lack the ability to comprehend the intrinsic dependencies within GraphQL. Instead, they focus on parameter variation and payload submission techniques [66] with the primary objective of identifying response errors (non-HTTP 200 status codes). In contrast, specialized tools like GraphCrawler [54], CrackQL [16] and GraphQL-Cop [28] aim primarily at testing individual queries and mutations without exploring the dependencies between them.

Putting dependency-aware testing aside, existing literature includes only a limited number of studies proposing methods to test GraphQL APIs. In white-box testing—which involves access to the API's underlying source code [30]—Zetterlund et al. [65] explored using production GraphQL queries to test the API. On the other hand, black-box testing—where testers lack access to the underlying source code [37]—has seen contributions by Vargas et al. [61] and Karlsson et al. [39], who proposed methods for inspecting the GraphQL schema and generating queries and mutations to test the API. A notable exception is Belhadi et al. [19], who introduce both black-box and white-box GraphQL fuzzers as plugins for the EvoMaster tool [20]. Their black-box approach involves randomly varying parameters of queries and mutations, while their white-box approach leverages source code information to generate queries
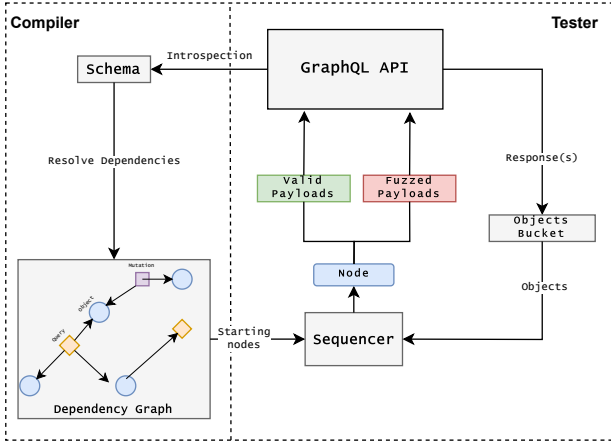
**Figure 2: Workflow of** GRAPHQLER**: The attacker's cycle represents the process of materializing a single node and sending it to the API.**

and mutations, aiming to produce a wide range of input scenarios to achieve comprehensive test coverage. Upon reviewing all existing literature on GraphQL testing, we found that EvoMaster's black-box testing is the only published tool that we could utilize and test against. Given their promising results, we employ it as one of the baseline methods for comparing the test coverage of our dependency-based approach with their parameter-varying methodology. It is essential to note that GRAPHQLER, as a black-box approach, cannot be directly compared with white-box methods.

Throughout our exploration of related work, it is evident that no tester in academic literature, open-source projects, or industry practices has leveraged dependency-based linking for testing GraphQL APIs. Instead, existing approaches mainly focus on modifying input parameters. Our proposed solution, GRAPHQLER, addresses this gap and elevates dependency-based testing to the forefront of automated GraphQL testing.

## 3 GRAPHQLER

This section introduces GRAPHQLER, our innovative approach to testing GraphQL APIs. This testing method is founded on a fundamental principle: the inherent dependencies among objects, mutations and queries in GraphQL. Instead of blindly sending requests to APIs, as is typical in traditional testing, GRAPHQLER leverages this dependency structure to intelligently chain queries and mutations relative to one another. By using a dependency-based testing approach, GRAPHQLER not only assesses response status codes (as conventional methods do) but also scans and stores actual responses, which then serve as inputs for subsequent requests.

As shown in Figure 2, GRAPHQLER begins by executing an *introspection query* to extract the schema of the target API. If the introspection query is unavailable, the schema is generated through brute-force using a word list [55]. This two-step process enables GRAPHQLER to operate as a *black-box* solution, allowing testing to proceed even without access to the introspection query. Once GRAPHQLER has the schema, the workflow unfolds in two distinct

phases. In the first phase, known as the *compilation phase*, the API schema is analyzed, and a dependency graph is constructed. The second phase, known as the *testing phase*, traverses this dependency graph using Depth-First Search (DFS) [58], materializes the payload, and sends requests to the server. During this phase, all created resources are tracked for later use in subsequent payloads.

### 3.1 Compilation Phase

The compilation phase is responsible for creating the dependency graph. The word "compilation" is used because it involves condensing the schema into meaningful sections. Given the extensive data provided by the schema, the focus of this phase is on extracting only the information essential for generating the dependency graph. The compilation phase involves three main steps:

**Annotating Mutations with Actions.** During the compilation phase, the schema is initially parsed to annotate all possible mutations with their corresponding actions. Unlike REST APIs, which use explicit HTTP methods (e.g., POST for creation, PUT for updating, DELETE for deletion), GraphQL lacks a built-in mechanism to indicate the intent of a mutation operation. Annotating mutations helps identify whether a mutation is intended for creating, updating, or deleting a resource. This identification is crucial for understanding the purpose of each mutation.

To identify each mutation, we examine both the method name and its documented description to identify action verbs for creation, updating, or deletion (e.g., createUser() would be annotated as a CREATE action, while destroyShip() with the description *"deletes a ship"* would be annotated as DELETE). If the action of a mutation cannot be determined, it is labeled as UNKNOWN.

**Inferring Dependencies.** In this step, we infer the relationships between different GraphQL types. While GraphQL explicitly defines dependencies between objects, it does not provide complete information about dependencies between objects and queries, or between objects and mutations. This complexity arises because both queries and mutations involve inputs and outputs, requiring dependencies to be considered on both ends.

The missing dependency to infer is the input to queries and mutations to the corresponding objects. To address this crucial link, we rely on GraphQL object IDs to link query and mutation inputs to their associated objects. This approach leverages the fact that GraphQL objects typically include a unique ID field. By examining the names and types of input fields, we can identify potential relationships to objects and document these dependencies.

It is essential to consider nullable types within GraphQL. When generating dependencies, if a type is marked as NON-NULL, we categorize it as a hard dependency, labeled as *hardDependsOn*; otherwise, it is a soft dependency, labeled as *softDependsOn*. This classification indicates the criticality of the dependency. Figure 3 provides a detailed illustration of the dependency resolution for each scenario, which we further explain below.

◇ *Object-Object Dependency.* For dependencies between objects, the resolution is straightforward, involving only a direct reference to the original schema. The *blue* line between User and Wallet in Figure 3 represents this type of dependency, showing that it is already defined in the schema and requires no additional
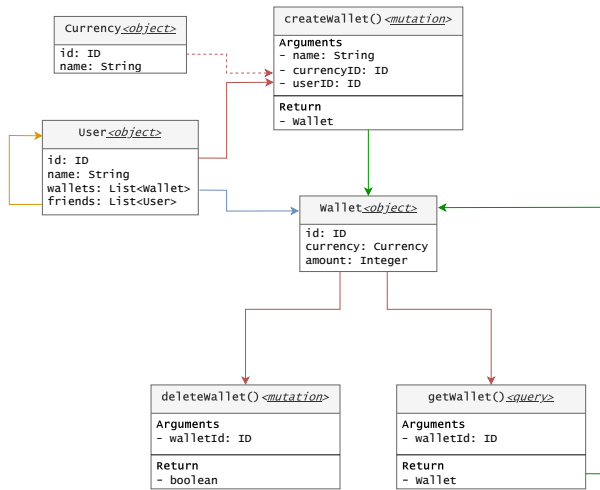
**Figure 3: Example of the types of dependencies in GraphQL.**

inference. Another form of object-object dependency occurs with self-references, depicted by the *orange* line. While this type also requires no additional inference, it does require special handling during dependency traversal.

◇ *Object-Mutation Dependency.* This dependency is divided into two categories: output dependencies and input dependencies. Output dependencies are established based on schema information. For input dependencies, we link the mutation to an object using either the mutation name or input fields such as ID. This is illustrated by the *red* arrows in Figure 3. In the same figure, the output dependency between the `createWallet` function and the `Wallet` object is marked in *green.* If no such link can be identified (e.g., when an `ID` field is absent), we classify it as an `UNKNOWN` dependency and make no connections.

◇ *Object-Query Dependency.* In this dependency type, both input and output dependencies are considered. The resolution process for these dependencies is similar to that used for object-mutation dependencies. This is shown in Figure 3 between the `getWallet` query and the `Wallet` object. Notably, this scenario forms a cyclical dependency, as the query only retrieves information about the object.

This structured approach ensures that the GraphQL compilation process thoroughly identifies and resolves dependencies. For a visual representation of annotating dependencies, refer to Figure 4, which illustrates a mutation's transition from a sample schema version to a version with resolved dependencies. In this example, NON_NULL marks the hard-dependency requirement. Two of the three fields are IDs, and they correlate to the `Currency` and `User` objects. Note that while the mutation has a hard dependency on the `Currency` object, the dependency on the `User` object is a soft dependency, as a `Wallet` can technically exist without an owner. The action associated with this mutation is determined to be CREATE, based on the mutation name, and the output is a `Wallet` object. After the annotations, we expect the dependency graph to have two directed edges linking the `User` and `Currency` objects
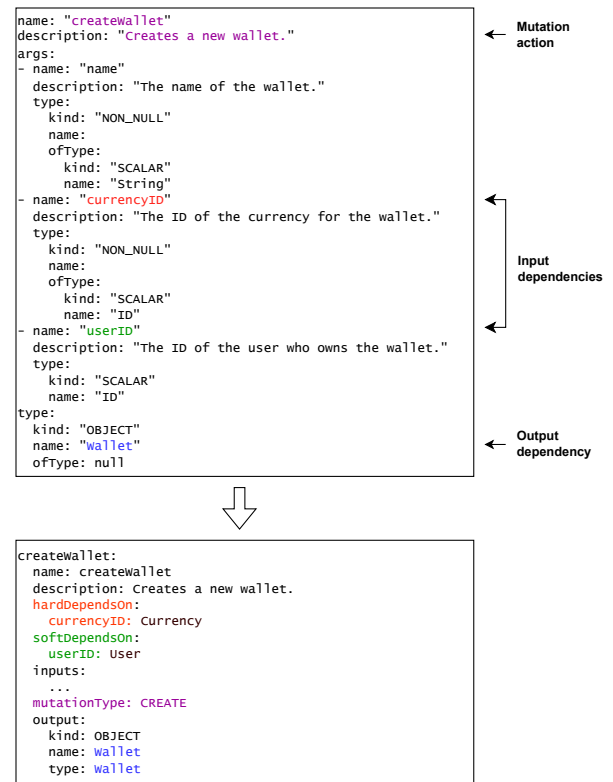


**Figure 4: Compilation of a sample `createWallet` mutation.**

to `createWallet`, and one edge from `createWallet` to the output object `Wallet`.

**Dependency Graph Creation.** Once dependencies are annotated on every query and mutation, we can begin constructing the dependency graph. This starts by generating a node for each query, mutation, and object. Each node has two main properties: the node's name (which corresponds to the name of the object/query/mutation) and the type (either object, query or mutation). For mutations, there is a third property, `mutation_type`, which denotes the action of the mutation as annotated in the previous steps.

With the nodes established, we create directed edges between them based on the compiled dependencies, where a parent node points to its dependent nodes. Note that, in this graph, connections between objects do not exist. This is because one must use a query or mutation to traverse the graph from one object node to the next. For example, to go from a `User` object to the `Wallet` object, one would need to use a new query called `getWallet()`, which accepts a `userID` as input and returns the wallet object as output. Note that an interesting case arises when objects are related to each other in the GraphQL schema. GraphQLer utilizes this fact during object storage and retrieval and uses objects to infer other objects.
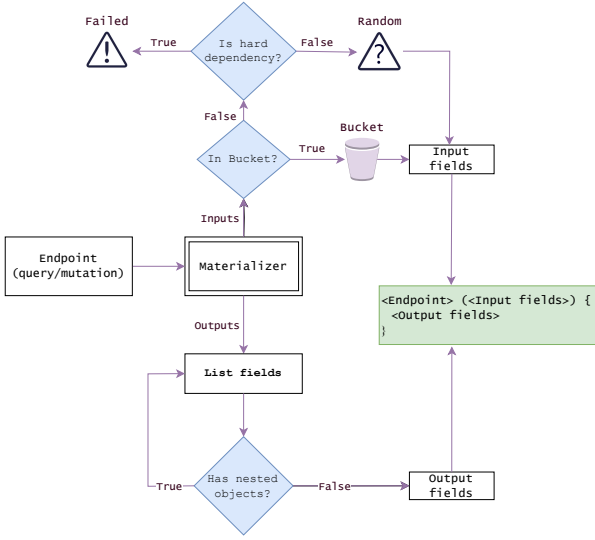
Omar Tsai, Jianing Li, Tsz Tung Cheung, Lejing Huang, Hao Zhu, Jianrui Xiao, Iman Sharafaldin, and Mohammad A. Tayebi



**Figure 5: The flow diagram depicting decision-making during the materialization step.**



**Figure 6: Payload generation for positive and negative testing.**

## 3.2 Testing Phase

The algorithm used by GRAPHQLER during the testing phase is detailed in Algorithm 1 of Appendix A. It begins by traversing the dependency graph and materializing the endpoint. This is then followed by sending payloads to the server and finally storing the returned resources for later use. Moreover, during this phase, GRAPHQLER tracks successful requests and responses—including any errors. If we consider the compilation phase as the "static" component, where each run produces the same output, then this phase represents the "dynamic" aspect of GRAPHQLER, where it actively interacts with the GraphQL API's responses, adapting its payload to accommodate the dynamic nature of the API's behavior. The testing phase is carried out in four main steps as follows:

**Initialization and Starting Points.** As the starting point of the testing phase, GRAPHQLER collects a list of independent nodes within the dependency graph, maximizing the number of entry points into the graph. This approach enhances the efficiency of resource reuse from the API and expands coverage. Notably, if no independent nodes are initially found—typically in fully cyclical graphs—GRAPHQLER seeks nodes with the fewest dependencies, starting from 1 (where 0 represents independent nodes). This iterative process continues until at least one starting node is identified, ensuring a well-rounded approach.

**Multi-Source Depth-First Search for Graph Traversal.** To navigate the dependency graph, GRAPHQLER uses a multi-source DFS algorithm (similar to the multi-source BFS [60])—starting with the multiple identified starting nodes. Our traversal algorithm considers the following aspects: 1) *Cycle Avoidance:* During traversal, we log visited nodes to avoid cycles and ensure efficient exploration of the graph; 2) *Traversal Monitoring:* We employ a stack to systematically track nodes during exploration, ensuring an organized approach to
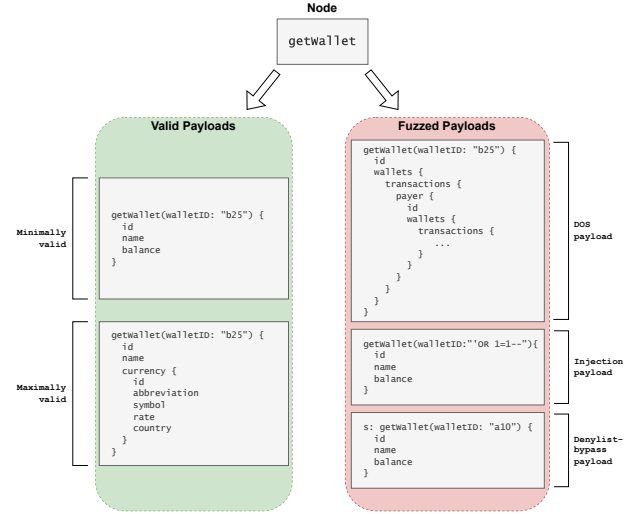
graph traversal; and 3) *Object Bucket Cache:* Concurrently, we maintain an object bucket cache to store any objects encountered during testing. This cache plays a pivotal role in preserving object references, ensuring a cohesive testing experience. The combination of these techniques streamlines the testing process, maintains order in graph traversal, and safeguards against cyclical dependencies.

**Materializing Requests.** When GRAPHQLER encounters a query or mutation node during traversal, it generates two types of payloads: valid and fuzzed. Each payload, as shown in Figure 6, is either a GraphQL query or mutation. The valid payloads include two variants: a *minimally valid* payload, containing the fewest output selectors to avoid overwhelming the API, and a *maximally valid* payload, which includes all selectors permitted by the endpoint. For fuzzed payloads, we explore several methods to test GraphQL endpoints for vulnerabilities. Figure 6 highlights three specific fuzzed payload examples—denial-of-service, SQL injection, and alias-based fuzzing—each carrying significant risks if a vulnerability is found. By systematically applying this A/B testing methodology, GRAPHQLER can compare standard API responses with those generated under adverse conditions, facilitating the identification of endpoint-specific vulnerabilities and deviations from expected behavior.

In the process of generating each payload, careful attention is paid to formatting both the input and output fields of the request. Within the input fields, two substitution scenarios are encountered. For scalar fields, GRAPHQLER performs a lookup in the bucket cache. If the specific scalar and its name have been previously recorded, GRAPHQLER replaces the scalar with a value that has been seen before. Conversely, if a lookup fails to yield a result, the field is filled with random data. For ID fields, object IDs are retrieved from the bucket cache; if an object's ID cannot be located, the materialization process fails, requiring the request to be deferred for later processing. It is also important to note that the nature of the random data generated for scalars depends on the data type
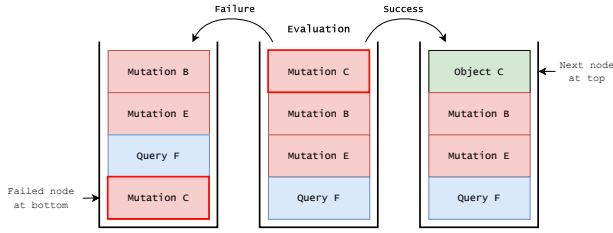
**Figure 7: Visualization of how the stack changes during traversal of the GraphQL endpoint.**

specified for the node. For instance, when the type is a string, a random string of arbitrary maximum length is produced; for booleans, randomized true or false values are generated, and this same logic applies to integers and floats.

Special consideration must be given to IDs to determine whether they represent a hard dependency. Materialization fails if an ID field is classified as a hard dependency and the corresponding ID is absent from the bucket. However, if the ID is not classified as a hard dependency, two possibilities arise: if the ID field is present in the bucket, it is substituted with the information stored therein; if not, the field is skipped, recognizing it as a non-critical dependency.

In output fields, one of GraphQL's beneficial characteristics is its ability to specify associated object details within the requested output. For instance, if we have a query named getUser() that returns a User object, and the User object has a dependent object called Role, we can specify both the attributes of the User and the Role in the output. GraphQLer uses this capability to recursively query for dependent objects of the query or mutation's output. This strategy allows us to retrieve as much information as possible, maximizing the coverage of the GraphQL API. Once the request is fully constructed, it is promptly dispatched to the GraphQL API for execution. Figure 5 provides a visual representation of both input and output materialization steps. Note that payload generation for both valid and fuzzed payloads follows a similar process. The key difference is that fuzzed payloads include malformed or malicious inputs based on specific attack types during their generation.

**Handling Server Responses.** When processing successful server responses, GraphQLer extracts and retains every unique object and scalar encountered and subsequently stores them in the object bucket cache. This process ensures the integrity of object references throughout the testing phase. When processing failed responses from the server, GraphQLer employs a dual approach. First, it attempts to address all failures returned, such as NON-NULL requirements on a payload's output field. This adaptive strategy aims to enhance the reliability and resilience of our testing process. Second, in cases where a response indicates an actual failure that cannot be immediately rectified, GraphQLer adopts a strategic approach: instead of prioritizing failed nodes at the top of the stack, it defers them to the bottom. This approach allows GraphQLer to focus on other nodes, preventing the testing process from becoming stuck due to recurring issues.

However, it is important to note that a reasonable limit on the number of retries is in place. Once this predefined threshold is

reached, the unsuccessful node is removed from further consideration and marked as failed. This controlled approach ensures that the testing process remains effective and efficient. Figure 7 provides a visual representation of how the stack is managed during DFS traversal. In the evaluation process of Mutation C, if the request is successful, the neighboring nodes of Object C will be added to the top of the stack. If the request fails, the original node, Mutation C, will be added to the bottom of the stack. The example highlights GraphQLer's ability to handle successes and failures. In practice, the traversal of the dependency graph is carried out in three distinct phases, each serving a specific purpose:

◇ *First Traversal.* During the initial phase, traversal focuses exclusively on CREATE mutations and regular queries. This selective approach helps maintain consistency within the objects bucket, ensuring that it is initialized in a controlled manner.
◇ *Second Traversal.* The second phase extends the scope of testing by introducing UPDATE mutations in addition to CREATE mutations and regular queries. This broadens the coverage of our testing and allows us to assess the behavior of more complex interactions.
◇ *Third Traversal.* In the final phase, all types of mutations and queries are encompassed, providing a comprehensive evaluation of the GraphQL API's capabilities and resilience.

As an integral part of the testing phase, GraphQLer systematically records and collects valuable data. This includes maintaining a list of queries and mutations that have been successfully executed, as well as capturing the responses. Additionally, GraphQLer maintains comprehensive logs of all queries and mutations sent to the server, ensuring a thorough record of the testing process.

## 4 Implementation

We have implemented GraphQLer in modularized Python, following best practices in software development. The compilation and testing phases are developed within the "compiler" and "fuzzer" modules, respectively. The critical link between these two modules is the dependency graph, which is generated at the conclusion of the compiler module and then loaded into memory at the start of the fuzzing module. Throughout the fuzzing phase, all inputs and outputs of the GraphQL API are logged for subsequent inspection. Additionally, the program provides real-time feedback by printing test results, which include the count of requests made, the number of successes, the number of failures, and any potential vulnerabilities found. For fine-grained control of how GraphQLer reacts to responses and for general settings, a settings file is provided in the repository, allowing users to fine-tune the tool to their liking. For visualization, GraphQLer generates a dependency graph to inspect API-identified dependencies, as shown in Figure 10 in the Appendix.

It is worth noting that in our implementation, we chose not to serialize GraphQL types into predefined Python types. This decision stems from the fact that the schema is already provided in dictionary form, and serialization would require subsequent deserialization during the materialization step of the testing phase—adding unnecessary complexity. Therefore, we opted to keep all data structures as Python dictionaries throughout the process.

**Table 1: Public GraphQL APIs used in baseline testing.**

| API | #Queries | #Mutations | #Objects |
|---|---|---|---|
| User Wallet [27] | 11 | 15 | 5 |
| Food Delivery [41] | 5 | 5 | 6 |
| Countries [11] | 6 | 0 | 5 |
| React-Finland [13] | 13 | 0 | 18 |
| Rick&Morty [14] | 9 | 0 | 7 |
| JSON-GraphQL [44] | 9 | 12 | 4 |
| GraphQLZero [17] | 13 | 19 | 18 |
| Anilist [18] | 27 | 29 | 122 |
| EHRI [38] | 19 | 0 | 46 |
| Universe [15] | 35 | 69 | 166 |
| PokeAPI [49] | 450 | 0 | 1952 |
| TCGDex [59] | 6 | 0 | 12 |

The implemented version of GraphQLer is a command-line tool that enables users to execute the compilation and testing phases either independently or sequentially through command-line flags. Running these phases separately provides a visualization of the dependency graph and allows for the manual annotation of unresolved dependencies.

## 5 Evaluation

In this section, we discuss the experimental setup and evaluation of GraphQLer. Specifically, we aim to address the following research questions:

◇ **RQ1:** Does GraphQLer efficiently cover possible GraphQL operations compared to baseline methods?
◇ **RQ2:** Can GraphQLer discover previously unknown GraphQL vulnerabilities in both open-source and commercial applications?
◇ **RQ3:** Does GraphQLer effectively explore the test space of a GraphQL operation in comparison to baseline methods, particularly in terms of execution time?

### 5.1 Experimental Design

*5.1.1 Evaluated Applications.* To conduct a comprehensive experimental evaluation, we test GraphQLer alongside the baseline methods using a diverse set of open-source and closed-source APIs, an intentionally vulnerable software and a commercial enterprise-level solution as outlined below.

◇ *Public APIs.* We utilize 12 public APIs across three categories: open-source, custom-built, and openly-hosted APIs [27, 34, 41]. For open-source APIs, we download the back-end code, set up a self-hosted GraphQL server and run tests. For custom-built APIs, we develop GraphQL APIs that mirror the structure of publicly accessible server-built APIs. For openly-hosted APIs, we reference free-to-use online APIs on GitHub and conduct direct testing.

◇ *Platform APIs.* This study evaluates five widely recognized proprietary GraphQL APIs: `Yelp` [10], `GitLab` [5], `Swop` [9], and `Blue Apron` [1].
We also conduct tests on a GraphQL API from a major financial institution, anonymized as `FinServ` under a non-disclosure agreement. This API supports millions of customers daily, providing real-time access to account information and transaction histories with minimal latency.

◇ *Vulnerable Software.* To assess GraphQLer's effectiveness, we use two applications with known GraphQL-related vulnerabilities to determine if GraphQLer successfully identifies these issues:
  ▷ `DVGA`. Damn Vulnerable GraphQL Application (`DVGA`) [4] is a deliberately insecure GraphQL API that provides a controlled environment for developers and security professionals to test and exploit vulnerabilities in GraphQL applications safely.
  ▷ `Saleor`. Saleor is an open-source, self-hosted platform that also offers a commercial version [8]. The `Saleor` API was found to be vulnerable to CVE-2022-39275 [3], which exposes broken access control through several GraphQL mutations. In testing the `Saleor` API, we aim to use previously reported CVEs associated with GraphQL APIs, focusing specifically on finding a reproducible CVE.

The diversity of our chosen systems under test allows us to validate GraphQLer's performance confidently.

*5.1.2 Baselines.* We compare our solution against the following baselines:

◇ GenericTester: This is a variation of GraphQLer in which no features are used—neither the dependency graph nor an object cache. In other words, the GenericTester just blindly tests queries and mutations, functioning like a generic API tester that uses random data to generate syntactically valid payloads.

◇ Evomaster: This method is the only published GraphQL testing approach for our comparative evaluation. It offers two approaches for assessing GraphQL APIs, with the relevant one being "black-box testing". Similar to GenericTester, EvoMaster [22] sends GraphQL mutation and query requests; however, it stands out by dynamically varying its payloads to explore a broader range of responses.

◇ ZAP: As a renowned open-source black-box scanner developed by OWASP, ZAP [47] serves as a robust tool for evaluating web vulnerability assessment and penetration testing. Although primarily for web applications, it also includes features for testing GraphQL APIs.

◇ BurpSuite: BurpSuite is a popular vulnerability scanner by PortSwigger [2]. To conduct GraphQL tests, we use the Auto GQL Scanner [31] extension as this is the only GraphQL extension that will automatically run payloads on GraphQL APIs.

*5.1.3 Evaluation Measures.* To address **RQ1**, we define a coverage metric for black-box testing that includes both positive and negative testing scenarios, as follows:

$$PositiveCoverage = \frac{\#NoErrors}{\#Endpoints}$$

$$NegativeCoverage = \frac{\#Errors}{\#Endpoints}$$

where *PositiveCoverage* represents the proportion of error-free endpoints that return data (beyond simply an HTTP 200 status), while *NegativeCoverage* indicates the proportion of endpoints that produce errors despite receiving valid inputs. These metrics address the success and failure conditions in black-box testing, where source code is inaccessible, making line coverage irrelevant. Additionally, we assess the diversity of response codes and returned objects in

**Table 2: Comparison of GraphQLer and baseline methods in coverage performance: The "+" and "-" signs indicate positive and negative coverage, respectively. A failure to execute the specific test is marked as FAILED.**

| API | GenericTester | | ZAP | | BurpSuite | | EvoMaster | | GraphQLer | |
|---|---|---|---|---|---|---|---|---|---|---|
| | (+) | (-) | (+) | (-) | (+) | (-) | (+) | (-) | (+) | (-) |
| User Wallet | 38.46% | 50.00% | 50.00% | 50.00% | 7.69% | 38.46% | 61.54% | 26.92% | **92.31%** | **100%** |
| Food delivery | 50.00% | 50.00% | 50.00% | 20.00% | 20.00% | 90.00% | **70.00%** | 40.00% | **70.00%** | **100%** |
| Countries | **50.00%** | 50.00% | 33.33% | 16.67% | 50.00% | **100%** | **50.00%** | 50.00% | **50.00%** | **100%** |
| React-Finland | **53.84%** | 46.15% | 16.67% | 23.08% | 0.00% | **100%** | 38.46% | 61.54% | **53.85%** | **100%** |
| Rick&Morty | **66.67%** | **100%** | 33.33% | **100%** | 0.00% | **100%** | **66.67%** | 33.33% | **66.67%** | **100%** |
| JSON-Graphql | 57.14% | 28.57% | 28.57% | 11.90% | 19.05% | 33.33% | 57.14% | 45.24% | **100%** | 85.71% |
| GraphQL Zero | 87.50% | 21.88% | 93.75% | 6.25% | 93.75% | 62.50% | 71.88% | 56.25% | **93.75%** | **93.75%** |
| Anilist | **8.92%** | 92.86% | FAILED | FAILED | FAILED | FAILED | FAILED | FAILED | **8.92%** | **94.64%** |
| EHRI | 42.10% | 84.21% | 10.53% | 36.84% | 0.00% | 31.58% | 84.21% | 57.89% | **94.74%** | **89.47%** |
| Universe | 19.23% | 91.30% | FAILED | FAILED | 0.00% | 12.32% | 14.01% | 88.46% | **82.24%** | **94.20%** |
| PokeAPI | **33.33%** | 85.84% | FAILED | FAILED | FAILED | FAILED | FAILED | FAILED | **33.33%** | **92.81%** |
| TCGDex | **100%** | 66.67% | 66.67% | 33.33% | 33.33% | **100%** | **100%** | **100%** | **100%** | **100%** |

the GraphQL API, which can enhance coverage and improve the likelihood of detecting bugs or vulnerabilities.

All tests are conducted on a server equipped with 16 CPU cores, 64 GB of RAM, and the Ubuntu 22.04 operating system, ensuring sufficient resources for consistent experimentation.

## 5.2 Test Coverage

The evaluated methods are tested on public APIs with varying queries, mutations and objects, as shown in Table 1. The results in Table 2 show that GraphQLer performs as well as or better than the baselines in both positive and negative coverage.

In terms of *PositiveCoverage*, GraphQLer outperforms EvoMaster in seven APIs and matches its performance in four. Compared to ZAP, GraphQLer exceeds its performance in all but two APIs. Similarly, GraphQLer consistently outperforms BurpSuite in terms of positive coverage, achieving better results across all evaluated APIs. Coverage was identical between GraphQLer and GenericTester for 6 APIs.

The performance gap between GraphQLer and all baselines, including ZAP, EvoMaster and, BurpSuite, becomes more pronounced with APIs featuring a greater number of mutations. This is because a higher number of mutations allows GraphQLer to track a broader range of resources, which can then be leveraged in subsequent server requests to generate more valid queries, thereby enhancing positive test coverage.

On average, GraphQLer achieves a 35% increase in *Positive-Coverage* over EvoMaster, with a maximum increase of over 60%. Compared to ZAP, GraphQLer shows a 38% average increase, with a maximum improvement of over 84%. Compared to BurpSuite, GraphQLer achieves an even higher average increase of 42%, with a maximum improvement exceeding 88%.

For the APIs that GraphQLer performed equally well with, such as the Countries API, Rick&Morty API, and the TCGDex API, the root cause is clear: These APIs contain a high number of queries that are unrelated to each other. Furthermore, they have highly specific input requirements that cannot be found from another query or mutation, necessitating users with additional knowledge to input the correct fields. A similar issue applies to the worst-performing API, Anilist, where search terms have to conform to specific show names, or else the API returns an HTTP 404 error.

To further emphasize the effectiveness of a dependency-based approach, the positive coverage increase from GraphQLer over GenericTester is 13.64%. Coverage increases were again the most significant for APIs with a larger number of mutations. This trend follows similar conclusions in the comparison with EvoMaster, ZAP, and BurpSuite.

For *NegativeCoverage*, GraphQLer consistently outperforms all baselines, as it is designed to generate both fuzzed and valid queries targeting GraphQL-specific vulnerabilities. We also observe that ZAP, EvoMaster, and BurpSuite achieve no coverage or fail on large APIs such as Anilist and PokeAPI, underscoring GraphQLer's strength in handling large schemas. The only instance where GraphQLer matches the negative coverage of ZAP is with the Rick & Morty API, which rejects any complex payloads, resulting in error responses. In conclusion, these results affirm that GraphQLer effectively explores the GraphQL API state space, addressing **RQ1**.

*5.2.1 Ablation study.* To validate that our solution—combining a dependency graph with an objects bucket—effectively improves testing coverage, we conducted an ablation study analyzing the

**Table 3: Potential vulnerabilities detected by GraphQLer and baseline methods in Public and Platform APIs (Q: Query, M: Mutation).**
✓ **: Detected;** ✗ **: Not detected.**

| API | Method | ZAP | BurpSuite | Evomaster | GraphQLer |
|---|---|---|---|---|---|
| React-Finland | Q: Conference | ✓ | ✗ | ✓ | ✓ |
| | Q: Contact | ✗ | ✗ | ✓ | ✓ |
| | Q: Interval | ✗ | ✗ | ✓ | ✓ |
| | Q: Schedule | ✗ | ✗ | ✓ | ✓ |
| | Q: Series | ✓ | ✗ | ✓ | ✓ |
| | Q: Theme | ✗ | ✗ | ✓ | ✓ |
| Anilist | Q: Thread | ✗ | ✗ | ✗ | ✓ |
| Gitlab | Q: runnerSetup | ✗ | ✓ | ✗ | ✓ |
| | Q: ciCatalogResources | ✓ | ✗ | ✗ | ✓ |
| | M: jiraImportStart | ✗ | ✗ | ✗ | ✓ |
| Yelp | Q: search | ✗ | ✗ | ✓ | ✓ |
| | Q: event_search | ✗ | ✗ | ✗ | ✓ |
| Blue Apron | Q: menu | ✗ | ✗ | ✗ | ✓ |
| | Q: subPauseDetails | ✗ | ✗ | ✗ | ✓ |

**Table 4: Comparison of GraphQLer and baseline methods for detecting vulnerabilities in DVGA.** ✓ **: Detected;** ✗ **: Not detected.**

| Vulnerability Category | ZAP | BurpSuite | Evomaster | GraphQLer |
|---|---|---|---|---|
| *Denial of Service* | | | | |
| Batch Query Attack | ✗ | ✗ | ✗ | ✓ |
| Deep Recursion Query Attack | ✗ | ✗ | ✗ | ✓ |
| Resource Intensive Query Attack | ✗ | ✗ | ✗ | ✓ |
| Field Duplication Attack | ✗ | ✗ | ✗ | ✓ |
| Aliases Based Attack | ✗ | ✗ | ✗ | ✓ |
| Circular Fragment | ✗ | ✗ | ✗ | ✓ |
| *Injection & Server-Side Vulnerabilities* | | | | |
| OS Command Injection | ✗ | ✓ | ✗ | ✓ |
| Server Side Request Forgery | ✗ | ✓ | ✗ | ✓ |
| SQL Injection | ✗ | ✓ | ✗ | ✓ |
| Path Traversal | ✗ | ✓ | ✗ | ✓ |
| *Client-Side & Other Vulnerabilities* | | | | |
| Information Disclosure | ✓ | ✓ | ✓ | ✓ |
| Stored XSS | ✗ | ✓ | ✗ | ✓ |
| HTML Injection | ✗ | ✓ | ✗ | ✓ |
| Query Deny List Bypass | ✗ | ✗ | ✗ | ✓ |

individual contributions of each component. Results from this analysis (see Table 5 in Appendix A) show that using both components together increases positive coverage by 27% over using only a dependency graph, and by 22% over using only the objects bucket. Negative coverage also increases by 30% and 19%, respectively, reflecting a more exhaustive exploration of the API. When the objects bucket is used without the dependency graph, improvements in coverage are essentially random: objects stored from previous responses are reused only by chance due to the lack of guided request sequencing. Conversely, without the objects bucket, the dependency graph alone is insufficient, resulting in behavior no better than sending isolated requests to each endpoint. By isolating each component of GraphQLer, we demonstrate that a dependency-based approach requires both state tracking and effective request ordering to comprehensively test GraphQL APIs.

## 5.3 Vulnerability Detection

Expanding API test coverage, particularly for complex interfaces like GraphQL, increases the likelihood of discovering vulnerabilities by testing a wider range of inputs, operations, and edge cases. Broader coverage helps reveal issues such as unhandled inputs, error-handling flaws, authorization gaps, and resource management weaknesses. This section presents GraphQLer's vulnerability detection results across various test scenarios, comparing its performance to baseline methods.

**Public & Platform APIs.** The detection of vulnerabilities in Public & Platform APIs is carried out through a detailed analysis of error responses during testing. Internal errors, particularly those manifesting as HTTP 500 responses, often point to underlying issues ranging from unhandled exceptions to severe security vulnerabilities, such as remote code execution. These error responses are critical indicators of flaws in the API's internal processing logic or infrastructure. Given the potential severity of these vulnerabilities,

such errors require immediate escalation and resolution by the development team.

In addition to HTTP 500 errors, it is also essential to investigate scenarios where the API returns HTTP 200 responses, but with embedded error messages or unexpected behavior. While these responses do not indicate a direct failure, they may still expose sensitive information or unintended behavior that could be exploited by an attacker. Consequently, such responses warrant a thorough review to identify security weaknesses, misconfiguration, or other risks that could compromise the API's integrity. This method of analyzing error responses offers a proactive approach to identifying vulnerabilities early in the development lifecycle, ensuring that potential security gaps are addressed promptly before they can be exploited.

Table 3 summarizes the potential vulnerabilities discovered by GraphQLer and other evaluated methods for public and platform APIs. In comparison with baseline methods, GraphQLer was able to identify errors that neither ZAP, EvoMaster, nor BurpSuite could find. Out of the 16 tested public and platform APIs, GraphQLer detects potential vulnerabilities in four of them, while ZAP and EvoMaster missed some. In total, GraphQLer detects 11 potential vulnerabilities, whereas ZAP and EvoMaster identify only three and seven vulnerabilities, respectively. Notably, BurpSuite was the worst of all, finding only a single error across the five APIs.

The effectiveness of GraphQLer in discovering potential vulnerabilities, compared to baseline methods, can be attributed to its unique approach of reusing information from returned responses. By leveraging this capability, GraphQLer generates more targeted payloads for fuzzing, resulting in greater overall coverage and an increased likelihood of vulnerability discovery. For example, in the case of the mutation error detected in GitLab, neither ZAP nor EvoMaster were able to identify it, while GraphQLer succeeded. This highlights GraphQLer's capability, particularly in APIs with
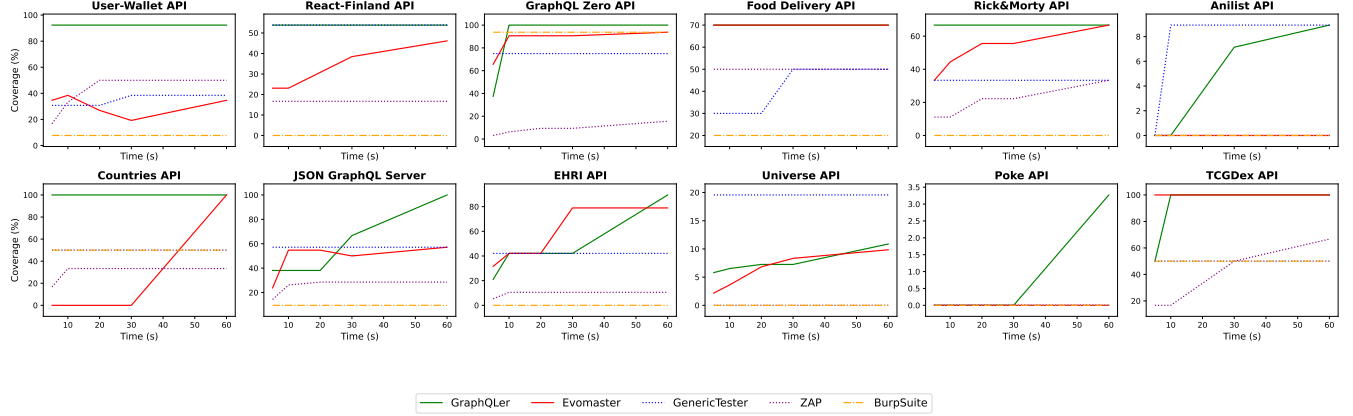
Figure 8: Positive coverage performance vs. time for GRAPHQLER and baseline methods.

mutations, where the sequence of mutations can be further refined based on their actions.

**DVGA.** As previously discussed, `DVGA` is an intentionally vulnerable application designed to include various classes of vulnerabilities. In our evaluation, we group these vulnerabilities into three main categories: *Denial of Service*, *Injection and Server-Side Vulnerabilities*, and *Client-Side and Other Vulnerabilities* (refer to Section 2). Table 4 summarizes the results of vulnerability detection across tools based on these categories. During testing, GRAPHQLER successfully identified all vulnerabilities present in DVGA. In contrast, both EvoMaster and ZAP were limited to detecting only the initial open introspection query vulnerability and failed to identify the rest. BurpSuite performed significantly better, uncovering nearly all common API vulnerabilities using its active scanner suite; however, it still missed several GraphQL-specific vulnerabilities. Note that, unlike EvoMaster, ZAP and BurpSuite—which depend on introspection query being enabled—GRAPHQLER operates independently of this limitation.

**FinServ.** For the `FinServ` API, we identified eight potential vulnerabilities, at least two of which were attacks leading to a denial of service. Each of these vulnerabilities exposed implementation details and sensitive information about the underlying API, either through stack traces or error codes. Moreover, GRAPHQLER's ability to run on the `FinServ` API without additional documentation demonstrates its capability to perform effective black-box fuzzing on large APIs.

**Saleor.** We self-host the `Saleor` API, re-deploying the commit hash that has the CVE along with all the necessary dependencies. GRAPHQLER quickly identified the CVE – a broken access control vulnerability. This includes all four vulnerable mutations along with the associated CVEs, generating the same error codes outlined in the advisory.

Based on these experiments, we conclude that GRAPHQLER outperforms other tools in identifying critical errors in the tested APIs, raising security concerns as these vulnerabilities could be exploited for various attacks. During testing, we encountered multiple instances where alternative tools crashed or terminated prematurely,

undermining reliability. For example, ZAP crashed while testing Yelp, and EvoMaster encountered a Java heap error with GitLab's API. In contrast, GRAPHQLER proved to be resilient, handling various GraphQL APIs, including both open-source and large-scale realworld applications. This confirms that GRAPHQLER can uncover previously unidentified vulnerabilities across diverse applications more effectively than the baselines, addressing **RQ2**.

## 5.4 Time Efficient Testing

Testing GraphQL APIs can be highly time-consuming due to their inherent size and complexity. A key factor in the effectiveness of a GraphQL API testing tool is its ability to efficiently cover a wide range of queries and mutations, which directly impacts the speed, scalability, and overall quality of the software development lifecycle. This is particularly important in CI/CD pipelines, where rapid feedback on code changes is critical to maintaining development velocity. Additionally, in security red-team engagements, where testing time is often limited, achieving broad test coverage within a constrained timeframe becomes essential. In such cases, optimizing for time-efficient testing enables comprehensive vulnerability detection and quality assurance without sacrificing coverage.

GRAPHQLER demonstrates exceptional time efficiency, consistently achieving higher coverage in much shorter durations—across time limits of 5 seconds to 1 minute—compared to baseline methods, as shown in Figure 8. One interesting observation is that GenericTester quickly converges to a fixed coverage percentage but does not progress beyond that point. This limitation is easily understood: as a generic tester aims to reach as many endpoints as possible without a strategic approach, it is primarily limited by network latency and I/O. Consequently, while GenericTester quickly achieves convergence, it ultimately gains a lower coverage score than other testers.

Across all tested APIs, we find that GRAPHQLER achieves higher coverage faster compared to ZAP, EvoMaster, and BurpSuite and higher eventual coverage compared to GenericTester. A key feature of GRAPHQLER is its use of the dependency graph to sequence tests, ensuring thorough coverage of queries and mutations while avoiding repetition, thus speeding up testing. In contrast, EvoMaster

and ZAP focus on generating diverse payloads without prioritizing their order. These findings highlight the effectiveness of our path traversal methodology, which accelerates testing and improves request quality. In conclusion, these results provide strong evidence for the effectiveness of our approach, fully addressing **RQ3**.

## 5.5 Diversity of Responses

An important measure of a GraphQL tester's performance is the diversity of objects it can retrieve. The more unique objects retrieved, the better, as this indicates that the tester is not substituting parameters based on pre-existing knowledge but is dynamically understanding the API and intelligently generating payloads. To assess the effectiveness of GraphQLer's context-aware fuzzing components, we compare the diversity of queried objects between GraphQLer and the baselines, which lack these advanced features. Across 12 public APIs tested, GraphQLer either retrieved the same or a greater number of unique objects than all testers, as shown in Figure 9. In four cases, none of the testers retrieved objects from the API, despite receiving HTTP 200 responses—likely due to the high level of domain knowledge required to craft effective payloads for these APIs. This result highlights the importance of the contextually aware testing approach implemented in GraphQLer enabling it to generate more diverse responses.

With context-aware fuzzing, research can be applied to diverse test cases, enabling analysis of how one request's outcome affects another. This helps uncover order-dependent failures like IDOR, a serious issue linked to broken access control. GraphQLer detects IDOR in GraphQL APIs by tracking cases where one user creates a private object that another can access. Its success in detecting CVE-2022-39275 and potential to identify authentication flaws demonstrates its effectiveness. Adding session management to GraphQLer enables effective testing for resource leaks, using techniques similar to REST API fuzzing [33].

Developing an efficient GraphQL tester requires addressing a key challenge: dependency chaining. This process inherently involves tracking resources across requests. GraphQL's complexity and interrelated components provide rich opportunities for deeper testing—yet, remarkably, no tools currently assess GraphQL APIs with this level of contextual awareness. To our knowledge, GraphQLer is the first tool to address this need for a context-aware approach, uniquely positioning itself to evaluate GraphQL APIs by mapping dependencies and caching objects effectively. As we observed, enabling a dependency-aware method allows for not only time-efficient testing strategies but also improvements in both errors discovered and the diversity of responses received. This enables GraphQLer to significantly enhance API test coverage and lay the groundwork for more comprehensive testing approaches in the future, as discussed below.

## 6 Discussion

**Internal Server Errors.** Internal server errors typically occur under specific conditions or actions. Some APIs respond with standard HTTP 200 statuses, as GraphQL allows an error section in the response body. An effective extension for investigating these internal server errors could involve integrating the source code with GraphQLer in a white-box approach to identify the root causes of
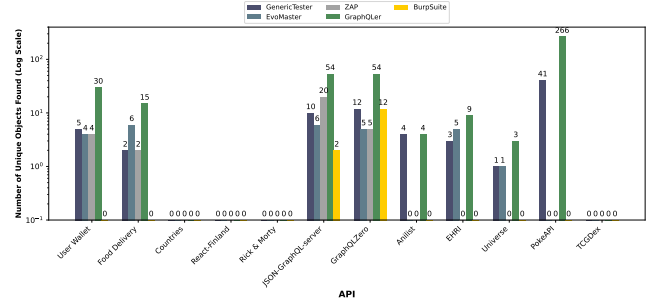


**Figure 9: Number of unique objects discovered by different testing tools.**

issues. Given that GraphQLer is a context-aware tester, leveraging the schema defined in back-end documentation would likely enhance the discovery of vulnerabilities.

**LLM-Driven Fuzzing.** Achieving a robust semantic understanding of complex domain-specific APIs remains a significant challenge. For example, APIs like `Anilist` require deep domain knowledge, where both GraphQLer and other tools show minimal coverage due to the uniqueness of each endpoint, making it hard to generate meaningful payloads. To address this, future work will explore leveraging large language models (LLMs) trained on diverse API patterns to improve GraphQLer's ability to generate semantically-aware queries. The ability of large language models to generate contextually valid payloads that can serve as an initial set of payloads can be used in conjunction with GraphQLer to create even more targeted attacks. The flexibility of GraphQL further complicates semantic fuzzing, but LLM-assisted techniques [32] provide a strong foundation. Key architectural components, such as inferring Object-Mutation and Object-Query dependencies, could be further refined through a deeper semantic understanding of the API using large language models. Additionally, advanced strategies like smart fuzzing, which use the object cache as context for a large language model, offer promising opportunities to uncover even more hidden vulnerabilities.

**End-to-end CI Solution.** GraphQLer can be easily integrated into any GraphQL-based continuous integration (CI) pipeline, requiring only the API endpoint to function. Since GraphQL allows flexibility between schema definitions and actual implementation, discrepancies between the schema and returned responses are not uncommon. GraphQLer is highly effective at detecting these inconsistencies, as well as identifying potential vulnerabilities that may arise from recent changes in the codebase.

**Ethical Considerations.** In our research, we prioritized ethical practices by ensuring that all identified vulnerabilities and potential bugs in the APIs were promptly and privately disclosed to the respective vendors. This approach aligns with responsible disclosure practices, providing vendors adequate time to address the issues and enhance the security of their APIs before any public dissemination of our findings. Our goal is to contribute positively to the cybersecurity community while minimizing potential risks to end users.

# 7 Conclusion

GraphQL testing is pivotal for ensuring the reliability, security, and performance of GraphQL APIs. Despite its critical significance, there is a significant gap in the availability of efficient testing tools in this domain. This paper addresses this gap by introducing GraphQLer, the first context-aware testing method for GraphQL APIs. By resolving dependencies and performing object caching, GraphQLer chains GraphQL payloads, ensuring that queries and mutations are executed with the expected inputs as originally intended. Our extensive experimental evaluation demonstrates the remarkable effectiveness of GraphQLer, surpassing the coverage performance of the best-performing baseline method by an average of 35%. Additionally, GraphQLer effectively identifies known CVEs and uncovers potential vulnerabilities in large-scale production APIs. Ultimately, the development of GraphQLer not only enhances the security posture of GraphQL APIs through robust automated testing but also lays the foundation for future advancements in GraphQL security tools for both scholars and practitioners.

# References

[1] Blueapron. https://hackerone.com/blueapron/policy_scopes. Accessed November 2, 2024.

[2] Burpsuite documentation. https://portswigger.net/burp/documentation.

[3] Cve-2022-39275. https://www.cve.org/CVERecord?id=CVE-2022-39275. Accessed June, 2024.

[4] Damn vulnerable graphql application. https://github.com/dolevf/Damn-Vulnerable-GraphQL-Application. Accessed November 2, 2024.

[5] Gitlab. https://www.gitlab.com. Accessed February 3, 2024.

[6] Graphql foundation. https://graphql.org/foundation/. Accessed October 15, 2023.

[7] The graphql report 2024. https://hygraph.com/resources/graphql-report-2024. Accessed November, 2024.

[8] Saleor. https://saleor.io/. Accessed February 3, 2024.

[9] Swop. https://swop.cx/. Accessed February 3, 2024.

[10] Yelp. https://www.yelp.com. Accessed February 3, 2024.

[11] Countries graphql api. https://countries.trevorblades.com/, 2023.

[12] Graphql voyager, 2023.

[13] react-finland. https://api.react-finland.fi/graphql, 2023.

[14] Rick&morty. https://rickandmortyapi.com/graphql, 2023.

[15] Universe api. https://www.universe.com/graphiql, 2023.

[16] Nick Aleks. Crackql - a graphql password brute-force and fuzzing utility. https://github.com/nicholasaleks/CrackQL. Accessed November 5, 2023.

[17] Emilio Almansi. Graphqlzero: Fake online graphql api for testing and prototyping, 2023.

[18] AniList. Anilist graphql api. https://github.com/AniList/ApiV2-GraphQL-Docs, 2023.

[19] Andrea Arcuri. EvoMaster: Evolutionary multi-context automated system test generation. In *2018 IEEE 11th International Conference on Software Testing, Verification and Validation (ICST)*. IEEE, apr 2018.

[20] Andrea Arcuri, Juan Pablo Galeotti, Bogdan Marculescu, and Man Zhang. Evomaster: A search-based system test generation tool. *Journal of Open Source Software*, 6(57):2153, 2021.

[21] V. Atlidakis, P. Godefroid, and M. Polishchuk. Restler: Stateful rest api fuzzing. In *Proceedings of the 41st International Conference on Software Engineering*, ICSE'19, pages 748–758, Piscataway, NJ, USA, 2019. IEEE Press.

[22] Asma Belhadi, Man Zhang, and Andrea Arcuri. Random testing and evolutionary testing for fuzzing graphql apis. *ACM Trans. Web*, August 2023.

[23] Fatna Belqasmi, Jagdeep Singh, Suhib Younis Bani Melhem, and Roch H. Glitho. Soap-based vs. restful web services: A case study for multimedia conferencing. *IEEE Internet Computing*, 16(4):54–63, 2012.

[24] K. Bennett, P. Layzell, D. Budgen, P. Brereton, L. Macaulay, and M. Munro. Service-based software: the future for flexible software. In *Proceedings Seventh Asia-Pacific Software Engeering Conference. APSEC 2000*, pages 214–221, 2000.

[25] Matteo Biagiola, Andrea Stocco, Filippo Ricca, and Paolo Tonella. Dependency-aware web test generation. In *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*, pages 175–185, 2020.

[26] Gleison Brito and Marco Tulio Valente. Rest vs graphql: A controlled experiment, 2020.

[27] Neo Chueng and Jianing Li. User wallet api. https://github.com/escobarflughafen/GraphQLer2022/tree/main/test_server. Accessed September 15, 2023.

[28] Dolev Farhi. Graphql cop - security audit utility for graphql. https://github.com/dolevf/graphql-cop. Accessed June 10, 2024.

[29] Roy Thomas Fielding and Richard N. Taylor. *Architectural Styles and the Design of Network-Based Software Architectures*. PhD thesis, University of California, Irvine, 2000. AAI9980887.

[30] Gordon Fraser, Matt Staats, Phil McMinn, Andrea Arcuri, and Frank Padberg. Does automated white-box test generation really help software testers? In *Proceedings of the 2013 International Symposium on Software Testing and Analysis*, ISSTA 2013, page 291–301, New York, NY, USA, 2013. Association for Computing Machinery.

[31] FWDSEC. Auto graphql scanner(auto gql). https://github.com/FWDSEC/burp-auto-gql/. Accessed November 8, 2023.

[32] Balaji Ganesan, Sambit Ghosh, Nitin Gupta, Manish Kesarwani, Sameep Mehta, and Renuka Sindhgatta. Llm-powered graphql generator for data retrieval. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 8657–8660, 2024.

[33] Patrice Godefroid, Bo-Yuan Huang, and Marina Polishchuk. Intelligent rest api data fuzzing. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ESEC/FSE 2020, page 725–736, New York, NY, USA, 2020. Association for Computing Machinery.

[34] Ivan Goncharov. Public graphql apis. https://github.com/graphql-kit/graphql-apis. Accessed September 15, 2023.

[35] Olaf Hartig and Jorge Pérez. Semantics and complexity of graphql. In *Proceedings of the 2018 World Wide Web Conference*, WWW '18, page 1155–1164, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee.

[36] Zac Hatfield-Dodds and Dmitry Dygalo. Deriving semantics-aware fuzzers from web api schemas. In *2022 IEEE/ACM 44th International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*, pages 345–346, 2022.

[37] Christopher Henard, Mike Papadakis, Mark Harman, Yue Jia, and Yves Le Traon. Comparing white-box and black-box test prioritization. In *Proceedings of the 38th International Conference on Software Engineering*, ICSE '16, page 523–534, New York, NY, USA, 2016. Association for Computing Machinery.

[38] European Holocaust Research Infrastructure. Ehri graphql api. https://portal.ehri-project.eu/api/graphql, 2023.

[39] Stefan Karlsson, Adnan Čaušević, and Daniel Sundmark. Automatic property-based testing of graphql apis, 2020.

[40] Leen Lambers, Lucas Sakizloglou, Osama Al-Wardi, and Taisiya Khakharova. Taint analysis for graph apis focusing on broken access control. In Russ Harmer and Jens Kosiol, editors, *Graph Transformation*, pages 180–200, Cham, 2024. Springer Nature Switzerland.

[41] Jianing Li. Sample food delivery graphql api, October 2023.

[42] Jiaxian Lin, Tianyu Li, Yang Chen, Guangsheng Wei, Jiadong Lin, Sen Zhang, and Hui Xu. forest: A tree-based approach for fuzzing restful apis, 2022.

[43] Qiang Liu, Flavio Toffalini, Yajin Zhou, and Mathias Payer. Videzzo: Dependency-aware virtual device fuzzing. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 3228–3245, 2023.

[44] Marmelab. json-graphql-server. https://github.com/marmelab/json-graphql-server, 2023.

[45] Phuong T. Nguyen, Juri Di Rocco, Claudio Di Sipio, Davide Di Ruscio, and Massimiliano Di Penta. Recommending api function calls and code snippets to support software development. *IEEE Transactions on Software Engineering*, 48(7):2417–2438, 2022.

[46] OWASP. Owasp top ten. https://owasp.org/www-project-top-ten/. Accessed February 3, 2024.

[47] OWASP. Zed attack proxy graphql. https://www.zaproxy.org/docs/desktop/addons/graphql-support/. Acessed November 8, 2023.

[48] Ina Papadhopulli and Hakik Paci. Graphql: A comprehensive analysis of its advantages, challenges, and best practices in modern api development. *The Eurasia Proceedings of Science Technology Engineering and Mathematics*, 32:230–237, 2024.

[49] Pokeapi. Pokemon api. https://studio.apollographql.com/public/poke-gql/variant/current/home, 2022.

[50] Antonio Quiña Mera, Pablo Fernandez, José María García, and Antonio Ruiz-Cortés. Graphql: A systematic mapping study. *ACM Comput. Surv.*, 55(10), February 2023.

[51] Antonio Quiña Mera, Pablo Fernandez, José María García, and Antonio Ruiz-Cortés. Graphql: A systematic mapping study. *ACM Comput. Surv.*, 55(10), feb 2023.

[52] Johannes Ryser and Martin Glinz. Using dependency charts to improve scenario-based testing. In *Proc. 17th Int'l Conf. Testing Computer Software*. Citeseer, 2000.

[53] Kenzaburo Saito, Yasuyuki Tahara, Akihiko Ohsuga, and Yuichi Sei. Proposal of an automated testing method for graphql apis using reinforcement learning. pages 1101–1107, 01 2025.

[54] Grant Smith. Graphcrawler. https://github.com/gsmith257-cyber/GraphCrawler. Accessed November 8, 2023.

[55] Nikita Stupin. Obtain graphql api schema even if the introspection is disabled. https://github.com/nikitastupin/clairvoyance. Accessed May 12, 2024.

[56] Ruben Taelman, Miel Vander Sande, and Ruben Verborgh. GraphQL-LD: Linked Data querying with GraphQL. In *Proceedings of the 17th International Semantic Web Conference: Posters and Demos*, October 2018.

[57] Sahar Tahvili, Marcus Ahlberg, Eric Fornander, Wasif Afzal, Mehrdad Saadatmand, Markus Bohlin, and Mahdi Sarabi. Functional dependency detection for integration test cases. In *2018 IEEE International Conference on Software Quality, Reliability and Security Companion (QRS-C)*, pages 207–214, 2018.

[58] Robert Tarjan. Depth-first search and linear graph algorithms. In *12th Annual Symposium on Switching and Automata Theory (swat 1971)*, pages 114–121, 1971.

[59] TCGdex. Tcgdex graphql api. https://www.tcgdex.dev/graphql. Accessed February 9, 2024.

[60] Manuel Then, Moritz Kaufmann, Fernando Chirigati, Tuan-Anh Hoang-Vu, Kien Pham, Alfons Kemper, Thomas Neumann, and Huy T. Vo. The more the merrier: efficient multi-source graph traversal. *Proc. VLDB Endow.*, 8(4):449–460, December 2014.

[61] Daniela Meneses Vargas, Alison Fernandez Blanco, Andreina Cota Vidaurre, Juan Pablo Sandoval Alcocer, Milton Mamani Torres, Alexandre Bergel, and Stéphane Ducasse. Deviation testing: A test case generation technique for graphql apis. In *Deviation Testing: A Test Case Generation Technique for GraphQL APIs*, 2018.

[62] Emanuele Viglianisi, Michael Dallago, and Mariano Ceccato. Resttestgen: Automated black-box testing of restful apis. In *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*, pages 142–152, 2020.

[63] Maximilian Vogel, Sebastian Weber, and Christian Zirpins. Experiences on migrating restful web services to graphql. In Lars Braubach, Juan M. Murillo, Nima Kaviani, Manuel Lama, Loli Burgueño, Naouel Moha, and Marc Oriol, editors, *Service-Oriented Computing – ICSOC 2017 Workshops*, pages 283–295, Cham, 2018. Springer International Publishing.

[64] S. V. Zelenov and S. A. Zelenova. Generation of positive and negative tests for parsers. *Program. Comput. Softw.*, 31(6):310–320, November 2005.

[65] Louise Zetterlund, Deepika Tiwari, Martin Monperrus, and Benoit Baudry. Harvesting production GraphQL queries to detect schema faults. In *2022 IEEE Conference on Software Testing, Verification and Validation (ICST)*. IEEE, apr 2022.

[66] Xiaogang Zhu, Sheng Wen, Seyit Camtepe, and Yang Xiang. Fuzzing: A survey for roadmap. *ACM Comput. Surv.*, 54(11s), September 2022.

## A  Appendix

### A.1  GRAPHQLER algorithm

### A.2  GRAPHQLER Coverage Ablation study

**Table 5: Coverage comparison across three GRAPHQLER variants. ODG stands for only dependency graph, whereas OOB stands for only objects bucket. The "+" and "−" indicate positive and negative coverage, respectively.**

| API | GRAPHQLER-ODG | | GRAPHQLER-OOB | | GRAPHQLER | |
|---|---|---|---|---|---|---|
| | (+) | (−) | (+) | (−) | (+) | (−) |
| User Wallet | 30.76% | 69.24% | 50.00% | 50.00% | 92.31% | 7.69% |
| Food delivery | 30.00% | 70.00% | 50.00% | 50.00% | 70.00% | 30.00% |
| Countries | 50.00% | 50.00% | 50.00% | 50.00% | 50.00% | 50.00% |
| React-Finland | 53.84% | 46.16% | 53.84% | 46.16% | 53.85% | 46.15% |
| Rick&Morty | 66.67% | 66.67% | 66.67% | 66.67% | 66.67% | 100% |
| JSON-Graphql | 57.14% | 28.57% | 61.90% | 28.57% | 100.00% | 85.71% |
| GraphQL Zero | 75.00% | 25.00% | 87.50% | 25.00% | 93.75% | 93.75% |
| Anilist | 8.92% | 94.64% | 8.92% | 94.64% | 8.92% | 94.64% |
| EHRI | 42.10% | 10.52% | 42.10% | 57.90% | 94.74% | 89.47% |
| Universe | 19.23% | 24.41% | 19.23% | 81.39% | 82.24% | 94.20% |
| PokeAPI | 33.33% | 3.70% | 33.33% | 66.67% | 33.33% | 92.81% |
| TCGDex | 50.00% | 33.33% | 50.00% | 66.66% | 100.00% | 100% |

### A.3  GRAPHQLER Artifacts

In the implementation of GRAPHQLER, one notable artifact generated is a graph that represents the dependencies among objects, mutations, and queries. Similar to [12], GRAPHQLER can visualize dependency graphs; however, it offers the added advantage of visualizing not only object dependencies but also the dependencies associated with mutations and queries.

---

**Algorithm 1** Main algorithm used in the testing phase of GRAPHQLER

**Require:** Starting nodes $StartingNodes$
**Ensure:** All reachable nodes are evaluated
1: Initialize $Visited \leftarrow \emptyset$
2: Initialize $ObjectsBucket \leftarrow$ empty mapping
3: Initialize $Stack \leftarrow StartingNodes$
4: **while** $Stack \neq \emptyset$ **do**
5:     $node \leftarrow Stack.\text{pop}()$
6:     **if** $node \in Visited$ **then**
7:         **continue**
8:     **end if**
9:     $(NextNodes, success) \leftarrow \text{Evaluate}(node)$
10:     **if** $success$ **then**
11:         $Stack \leftarrow Stack + NextNodes$
12:         $Visited \leftarrow Visited \cup \{node\}$
13:     **else**
14:         $Stack \leftarrow NextNodes + Stack$
15:         $\text{LogError}(node)$
16:     **end if**
17: **end while**
18: **function** EVALUATE($node$)
19:     **if** not dependenciesMet($node$) **then**
20:         **return** $(\{node\}, \textbf{False})$
21:     **end if**
22:     $payload \leftarrow \text{Materialize}(node)$
23:     $(success, data) \leftarrow \text{SendRequest}(payload)$
24:     **if** $success$ **then**
25:         PutObjectsInBucket($data$.objects)
26:         **return** $(node.\text{children}, \textbf{True})$
27:     **else**
28:         **return** $(\{node\}, \textbf{False})$
29:     **end if**
30: **end function**
31: **function** MATERIALIZE($node$)
32:     $inputs \leftarrow \text{MaterializeInputs}(node, ObjectsBucket)$
33:     $output \leftarrow \text{MaterializeOutput}(node, ObjectsBucket)$
34:     $payload \leftarrow inputs + output$
35:     **return** $payload$
36: **end function**
37: **function** DEPENDENCIESMET($node$)
38:     **for all** $dependency \in node.\text{dependsOn}$ **do**
39:         **if** $dependency \notin ObjectsBucket$ **then**
40:             **return** False
41:         **end if**
42:     **end for**
43:     **return** True
44: **end function**
45: **function** PUTOBJECTSINBUCKET($objects$)
46:     **for all** $object \in objects$ **do**
47:         $ObjectsBucket[object.\text{name}] \leftarrow ObjectsBucket[object.\text{name}] \cup \{object.\text{id}\}$
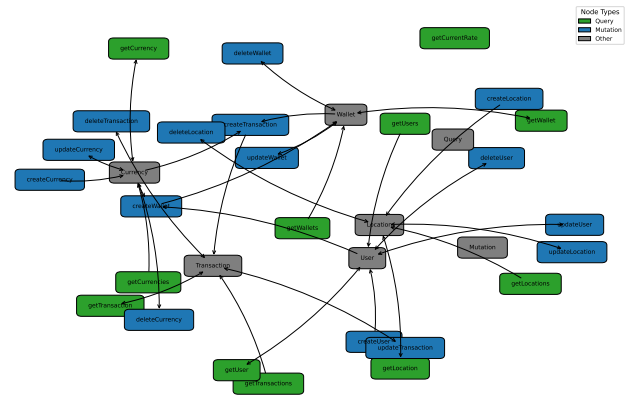48:     **end for**
49: **end function**



**Figure 10: Dependency graph generated by GRAPHQLER for the User Wallet API.**