# COTS: Connected OpenAPI Test Synthesis for RESTful Applications

Christian Bartolo Burlò[1]([✉]) [iD], Adrian Francalanza[2] [iD], Alceste Scalas[3] [iD], and Emilio Tuosto[1] [iD]

[1] Gran Sasso Science Institute, L'Aquila, Italy
christian.bartolo@gssi.it
[2] Department of Computer Science, University of Malta, Msida, Malta
[3] DTU Compute, Technical University of Denmark, Kongens Lyngby, Denmark

**Abstract.** We present a novel model-driven approach for testing RESTful applications. We introduce a (*i*) domain-specific language for OpenAPI specifications and (*ii*) a tool to support our methodology. Our DSL, called COPENAPI, is inspired by session types and enables the modelling of communication protocols between a REST client and server. Our tool, dubbed COTS, generates (randomised) model-based test executions and reports software defects. We evaluate the effectiveness of our approach by applying it to test several open source applications. Our findings indicate that our methodology can identify nuanced defects in REST APIs and achieve comparable or superior code coverage when compared to much larger handcrafted test suites.

## 1 Introduction

Modern software is increasingly composed of concurrent and distributed components that are independently developed, and function by exchanging data across a communication network. The interaction of such components may be based on classic client-server Internet protocols (such as SMTP, IMAP, and POP3), various forms of remote procedure calls (RPCs), web-based standards—such as the REpresentational State Transfer (REST) architectural design.
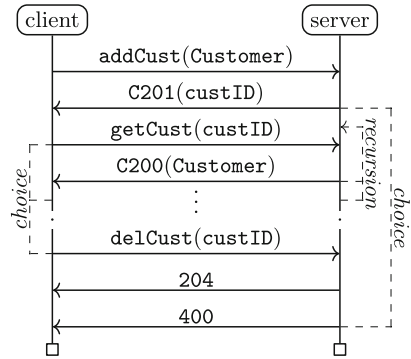
Ensuring the correctness and reliability of these applications is notoriously hard. Software developers typically create handcrafted test suites, which must address many (and ever-growing) application usage scenarios: It is often the case that such manually-written test suites are developed intermittently, over long periods of time, by a variety of testers. This makes the software development and testing process error-prone and susceptible to inconsistencies [23,33]. A number of efforts have emerged with the aim of simplifying the testing process of component based-software by streamlining the amount of manual work through automatic test generation. Most of these efforts focus on automatic test generation for applications exposing an API over a network, where there is a considerable interest in addressing the popular REST API style [4,5,7,13,15,21,30–33,37,43].

```
1   paths:
2     /customer:
3       POST:
4         operationId: addCust
5         responses: '201' ...
6     /customer/id:
7       GET:
8         operationId: getCust
9         responses: '200' ...
10      DELETE:
11        operationId: deleteCust
12        responses: '204' ...
13    /card:
14      POST:
15        operationId: addCard
16        responses: '201' ...
17    /address:
18      POST:
19        operationId: addAddr
20        responses: '201' ...
```



(a) OpenAPI spec. for an online store.

(b) Message Sequence Chart outlining the invocation dependencies in Fig. 1a.

**Fig. 1.** Typical documentation provided for REST APIs.

*Example 1.* Figure 1a shows a fragment of a REST API of the *SockShop* web application [38] written as an OpenAPI specification [28]. A new customer can be created by passing the necessary credentials as payload to the operation with ID `addCust` (lines 3–4); this operation returns a unique reference identifier for the customer, meant to be used in subsequent operations to create a payment card (`addCard`) and an address (`addAddr`) to associate with (previously created) customers. The customer information including any cards and addresses linked to it can be retrieved with `getCust` and deleted with `deleteCust`.                    □

The OpenAPI in Fig. 1a is a typical instance of a REST API description. OpenAPI is used pervasively for such specifications: it provides information on the available URIs, the corresponding HTTP methods and respective responses with the associated data formats (omitted from Fig. 1a). However, OpenAPI does not specify the relationships and dependencies between different API invocations. For instance, Fig. 1a does *not* express that:

– Customer information can be added (`addCard`, `addAddr`), retrieved (`getCust`), or deleted (`deleteCust`) only after the customer is created (`addCust`).
– Operations `addCard`, `addAddr` or `getCust` cannot be executed after the customer has been deleted (`deleteCust`).
– Operations `addCard`, `addAddr` and `getCust` may be interleaved without effecting the successful outcome of the other invocations.
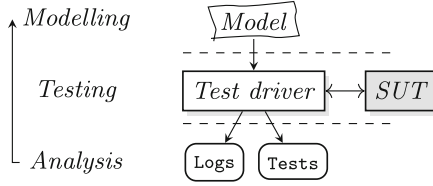
– The variable `custId`, an ID given to the customer upon creation is to be used to perform any operation on the customer (such as retrieval or deletion).

These dependencies induce an ordering in the API invocations, depicted in the message sequence chart in Fig. 1b. Unfortunately, this information is often omitted (cf. [8,19,38]) or only informally stated (as e.g., in [14]). In development contexts, this is often sufficient as developers are typically able to deduce the intended order of invocation by examining the descriptions of requests and the data types exchanged by each operation. However, this omission becomes problematic when testing. Deviations from the dependencies identified earlier in Fig. 1b can be categorized as *logic-based faults*, that is, errors stemming from incorrect or unintended actions within the business logic of the application. Such violations are often intricate and context-specific, arising from the unique interactions and workflows within the application. These faults are not easily detectable through basic system responses or standard error codes and typically necessitate comprehensive, scenario-based testing for identification. Furthermore, these faults are intrinsically linked to the essential operations and branching mechanisms of the application, impacting its ability to perform as intended.

Previous work on REST API testing proposes various *fully-automated* test generation strategies [5,21,33,42]. These *push-button* tools excel at detecting *systemic errors* that are related to the system's general functioning such as server crashes, malformed requests, unauthorized access attempts, or resource not found errors. However, it is unlikely that logic-based errors are uncovered by invoking the API without an explicit pre-defined model of the intended interaction with the API. *E.g.,* in the *SockShop* application from Example 1, such tools are able to detect if the SUT crashes when creating a new customer, but it is unlikely that they are able to detect whether the SUT allows the retrieval of a previously deleted customer. A test to detect such errors should create a customer and then try to retrieve the same customer after having deleted it. It has been empirically shown that the quality of tests substantially improves when conducted in a *stateful* manner, i.e., by exploring states of the system under test (SUT) that are reachable via sequences of invocations [5]. A recent work concludes that: *"existing tools fail to achieve high code coverage due to limitations of the approaches they use for generating parameter values and detecting operation dependencies"* [22].

The state-of-the-art of testing REST APIs can be seen on a spectrum. At one end, manual tests, though labor-intensive, excel at uncovering complex, logic-based errors. At the opposite end are fully-automated testing tools, which, while requiring minimal effort, fall short in detecting eloaborate errors. This work seeks a middle ground, leveraging a *model-based* approach. Our goal is to automate the identification of more logic-based faults, thus bridging the gap between labor-intensive manual testing and the scope of fully-automated tools.

*Contributions.* We present a model-based, tool-supported methodology for the automatic testing of web applications exposing REST APIs. Our approach allows the specification of API dependencies and constraints. Our contributions are:

**Fig. 2.** Methodology overview

1. COPENAPI: a domain-specific language designed to specify dependencies between requests in an OpenAPI specification that capture the state of the interaction with the SUT and the sequencing of message exchanges;
2. COTS: an automated tool leveraging COPENAPI models to generate tests that interact with the SUT and assess the correctness of its responses.
3. Experimental results demonstrating the effectiveness of COPENAPI and COTS using case studies from various application domains. We compare our results against RESTful API testing tools and manually-written tests.

The key benefits of our approach are in terms of (*i*) a high degree of expressiveness due to the possibility of specifying data dependencies, (*ii*) effectiveness of the approach that can identify logic-based faults, and (*iii*) high level of coverage.

## 2    A Model-Driven, DSL-Based Methodology

### 2.1    Background

Based on the REST architectural style [16] and HTTP protocol, REST APIs facilitate software services on the web by exposing *resources* through HTTP URIs, which are manipulated using standard HTTP *request methods* like GET, POST, PUT, and DELETE. Responses are represented by standard HTTP codes (e.g., 200 for success, 404 for missing resources, 500 for server errors). While OpenAPI and GraphQL [18] are prominent for documenting REST APIs, our emphasis on OpenAPI stems from its broad adoption and tool support for generating client libraries across various programming languages. Though our methodology primarily targets OpenAPI, it's flexible for other web API standards like GraphQL.

The technical nature of REST APIs, their complex state-dependent interactions, and their evolving nature introduces challenges for which existing model-based testing methods and tools are ill-suited (see related work in Sect. 5). To bridge this gap, we introduce the COPENAPI DSL (Sect. 3), designed to articulate the intricate sequences and state dependencies characteristic of RESTful interactions, enabling precise testing beyond the capabilities of existing models and tools.

## 2.2    Methodology

Our methodology consists of the three phases depicted in Fig. 2:

**(P1)** *modelling phase:* construct a model that describes how a client application
might interact with the SUT, and how the SUT is expected to respond to
the client's inputs;
**(P2)** *testing phase:* automatically generate a *test driver* that tests the SUT by
interacting with it according to the model, and reports whether the SUT
responses violate the model;
**(P3)** *analysis phase:* inspect the outputs produced by test driver to identify
faults in the SUT.

The phases **(P1)**, **(P2)** and **(P3)** are iterated when the model has to be refined
according to the findings of the analysis or it has to be evolved according to the
life cycle of the SUT. The main goal of our approach is to automate the testing
phase, and let testers shift their efforts to the modelling phase, reducing the need
to develop (and maintain) a large suite of handcrafted test.

The first step towards the methodology entails determining a suitable model
for phase **(P1)**, allowing for the automatic derivation of tests for RESTful appli-
cations. To write such model, we design COPENAPI (see Sect. 3) taking inspi-
ration from *(binary) session types* [20].

The next step in our methodology is **(P2)**, i.e., the testing phase. In this
phase, the test driver interacts with the SUT by sending and receiving messages
according to the COPENAPI model. The test driver needs to perform two types
of interactions with the SUT:

– the test driver must send requests to the SUT of the correct format and
  payload type; and
– the test driver needs to check that the responses received from the SUT do
  not violate the COPENAPI model. These violations may be of three forms:
    • the response code received is not valid, or
    • the payload data does not have the expected format, or
    • the response code and payload data format are valid, but the payload data
      causes an assertion violation in the model (e.g., by violating a constraint
      involving data from a previous request/response).

The last phase of the methodology Item **(P3)** is devoted to the identification
of faults in the SUT, by inspecting the outputs of the test driver. The test driver
produces a successful test when it manages to complete a full traversal of the
COPENAPI model without finding any of the aforementioned errors. To help
with the analysis, the test driver produces two outputs:

1. a log file with information about every performed test: the random seed used
   to generate it, the sequence of messages sent and received and their respective
   payloads, and the test outcome (pass/fail); and
2. an offline representation of failed tests, usable for manually reproducing faults
   of the SUT without re-executing the test driver, useful for, e.g., bug reporting.

## 3   A DSL for OpenAPI

We introduce a DSL based on an augmented form of *session types* [12,20,36].
Intuitively, a session type specifies the valid sequences of message exchanges (i.e.,
the *protocol*) that regulate the interaction between communication programs;
here we focus on *binary session types*, which model how a program is expected
to interact with just one other program (in our case, a client with a RESTful
server).

The formal syntax of our augmented session types $S$ is:

$$
\begin{aligned}
S ::= \; & +\left\{!operationId_i(\mathsf{gen}_i).S_i\right\}_{i\in I} && \textit{(internal choice)} \\
| \; & \&\left\{?responseCode_i(\mathsf{dec}_i)\langle A_i\rangle.S_i\right\}_{i\in I} && \textit{(external choice)} \\
| \; & \mathsf{rec}\ X.S && \textit{(recursion)} \\
| \; & X && \textit{(recursion variable)} \\
| \; & \mathsf{end} && \textit{(termination)} \\
T ::= \; & \mathtt{Int}\ |\ \mathtt{String}\ |\ \dots && \textit{(data types)} \\
G ::= \; & \mathtt{Int}(g)\ |\ \mathtt{String}(g)\ |\ \dots && \textit{(data types with generators)}
\end{aligned}
$$

where $I$ is a finite non-empty set of indexes and for $i \in I$, $\mathsf{gen}_i$ is a *generator*
assignment $x_{i,1} : G_{i,1}, \dots, x_{i,n_i} : G_{i,n_i}$ (cf. Example 2) and $\mathsf{dec}$ is a type assign-
ment $x_{i,1} : T_{i,1}, \dots, x_{i,n_i} : T_{i,n_i}$) with $x_{i,1}, \dots, x_{i,n_i}$ pairwise distinct. We use
standard data types, including standard types (e.g., $\mathtt{Int}$, $\mathtt{String}$, ...) and user-
defined types. Generators are also user-defined, tailored to the particular request
being sent. This enables crafting input scenarios that are specifically aimed at
exploring under-tested paths or complex behaviour of the API, enhancing both
the depth and breadth of the testing coverage.

The communication units of our syntax are the *output* and *input* prefixes,
respectively !*operationId* and ?*responseCode* where *operationId* is the corre-
sponding identifier referring to the specific request in the OpenAPI specification
and *responseCode* is the HTTP response code that is returned from the service
after the request is made (see Sect. 5). Assertions $A$ are boolean expressions that
may refer to the variables occurring in preceding payload descriptions.

We illustrate the usage of COPENAPI with the following example which
elaborates on the protocol in Fig. 1b.

*Example 2.* Figure 3 shows a testing model written in COPENAPI. The model
corresponds to the message-sequence diagram from Fig. 1b. The model specifies
that the interaction protocol should first create a customer by invoking `addCust`
together with two payload generators:

1. `genApiKey` which retrieves a key to authenticate with the API; and
2. `genCustInfo` to generate the customer information.

The server is then expected to reply with the response code `201` indicating
that the customer creation was successful, while also including an identification
code of type `String` which is bound with the variable `custId`. Next, the driver

```
1  S_shop = !addCust(apiKey: String(genApiKey), c1: Customer(genCustInfo)).
2           ?C201(custId: String).
3           rec X.(
4              +{ !addCard(apiKey, custId, card: Card(getCardInfo)) .
5                    ?C201(CardId: String) . X,
6                 !addAddr(apiKey, custId, addr: Address(getAddrInfo)) .
7                    ?C201(addressId: String) . X,
8                 !getCust(apiKey, custId) .
9                    ?C200(c2: Customer)<checkCustomer(c1,c2)> . X,
10                !deleteCust(apiKey, custId).?C204() . end
11             }
12          )
```

**Fig. 3.** Example of a COPENAPI model using the OpenAPI specification in Fig. 1a, formalising the message-sequence diagram in Fig. 1b.

chooses between sending `addCard`, `addAddr`, `getCust` or `deleteCust`, including the data required in the respective payloads by either invoking the specified generators, or use previously-bound variables. For instance, in the case where `addCard` is selected, the payload must include: *(1)* the API key stored in the variable `apiKey`, *(2)* the specific customer ID stored in the variable `custId`, and *(3)* the card information, by invoking the generator `getCardInfo`.

The model also specifies the use of a user-provided assertion `checkCustomer` (line 9) to check the response sent by the SUT for the response `getCust`. The assertion checks whether the retrieved customer information (stored in $c2$) matches with the information provided when the customer was created when `addCust` was invoked (stored in $c1$). If the assertion succeeds, the test driver loops and performs the choice (line 3) once again. Otherwise, the test run terminates and the test is marked as failed. The recursion repeats until the test driver selects `delCust`, after which the protocol would terminate successfully.     □

### 3.1 Semantics of COPENAPI

We give a denotational semantics of COPENAPI by mapping its model to sets of finite sequences of input/output *events*. Given a output prefix $\pi$, $[\![\pi@i]\!]$ denotes the set of values in the codomain of the generator of the $i$-th parameter of $\pi$ while for an input prefix $\pi$, $[\![\pi@i]\!]$ is the set of values inhabiting the type of the $i$-th parameter of $\pi$. We define

$$E = \{\pi(v_1, \ldots, v_n) \mid \pi \text{ is a prefix and } \forall 1 \le i \le n : v_i \in [\![\pi@i]\!]\}$$

Let $M$ a COPENAPI model and $\rho$ be a map assigning subsets of $E^\star$ to free recursion variables of $M$. The semantics of $M$ in $\rho$, is the set $[\![M, \rho]\!]$ defined according to the equations in Fig. 4. The semantics is defined by induction on the structure of the COPENAPI models and returns a set of finite traces; the idea is that this set represents all the executions allowed by a COPENAPI specification. Basically, $[\![M, \rho]\!]$ yields the possible tests that the test driver might use. For

for all $i \in I$, $\pi_i$ output prefix, $\boldsymbol{x}_i = x_{i,1}, \ldots, x_{i,n_i}$ and $\boldsymbol{v}_i = v_{i,1}, \ldots, v_{i,n_i}$

$$[\![ + \{ \pi_i(x_{i,1} : \mathsf{G}_{i,1}, \ldots, x_{i,n_i} : \mathsf{G}_{i,n_i}).\mathsf{S}_i \}_{i \in I} ]\!] \triangleq$$

$$\bigcup_{i \in I} \{ !\pi_i(\boldsymbol{v}_i).r_i \mid \forall 1 \leq j \leq n_i : v_{i,j} \in [\![ \pi_i @ j ]\!] : r_i \in [\![ \mathsf{S}_i[\boldsymbol{v}_i / \boldsymbol{x}_i], \rho ]\!] \}$$

for all $i \in I$, $\pi_i$ input prefix, $\boldsymbol{x}_i = x_{i,1}, \ldots, x_{i,n_i}$ and $\boldsymbol{v}_i = v_{i,1}, \ldots, v_{i,n_i}$

$$[\![ \& \{ \pi_i(x_{i,1} : \mathsf{G}_{i,1}, \ldots, x_{i,n_i} : \mathsf{G}_{i,n_i}).\mathsf{S}_i \}_{i \in I} ]\!] \triangleq$$

$$\bigcup_{i \in I} \{ \pi_i(\boldsymbol{v}_i).r_i \mid \forall 1 \leq j \leq n_i : v_{i,j} \in [\![ \pi_i @ j ]\!], A_i[\boldsymbol{v}_i / \boldsymbol{x}_i] \Downarrow \mathsf{tt}, r_i \in [\![ \mathsf{S}_i[\boldsymbol{v}_i / \boldsymbol{x}_i], \rho ]\!] \}$$

$$[\![ \mathsf{rec}\ X.\mathsf{S}, \rho ]\!] \triangleq \bigcap \{ R \mid [\![ \mathsf{S}, \rho[X \mapsto R] \subseteq R ]\!] \} \qquad [\![ X, \rho ]\!] \triangleq \rho(X) \qquad [\![ \mathsf{end}, \rho ]\!] \triangleq \{ \epsilon \}$$

**Fig. 4.** Semantics for the COPENAPI language.

example, the internal choice is defined as the set of all possible output prefixes in the choice with the different variations of the values that can be generated from the generators. The case for the external choice is similar, with the only exception that once the values are replaced in assertions of branches should hold. The semantics of end is standard. The map $\rho$ keeps track of the recursive variables in $\mathsf{S}$; initially we assume that $M$ is *closed*, i.e., it has no free occurrences of recursion variables. This makes the semantics of recursion standard.

### 3.2 Implementation

We implement the COPENAPI language as part of a tool called COTS. Implemented in the Scala programming language, takes as input the COPENAPI model and the optional preamble, and generates the Scala source code of an executable *test driver* that interacts with the SUT according to the model. The test driver, in turn, interacts with the REST API exposed by the SUT by using a Scala API which is autogenerated from the provided OpenAPI specification, using OpenAPI Generator.[1] When the test driver runs, it invokes such Scala API methods to send HTTP requests to the SUT, and to receive and parse its responses; the model determines which requests are sent (and in what order) by the test driver, and which responses are expected. After completing the test runs, the test driver produces the following output:

*(1)* the level of model coverage achieved by the test runs;
*(2)* a log file with information about every performed test: the random seed used to generate the test, the sequence of requests/responses and their payloads, and the test outcome (pass/fail);
*(3)* an offline representation of failed tests as sequences of `curl`[2] commands, which can be executed from a shell to reproduce faults of the SUT without re-executing the test driver.

---

[1] https://openapi-generator.tech.
[2] https://curl.se.

**Table 1.** Case studies.

| Application | Notes |
|---|---|
| *FeaturesService* [14] | Used in empirical evaluation of [4, 24, 33, 41] |
| *RESTCountries* [27] | Used in the empirical evaluation of [41] |
| *GestaoHospital* [35] | Used in the empirical evaluation of [41] |
| *LanguageTool* [29] | Used in the empirical evaluation of [24, 41] |
| *PetClinic* [8] | Used in the empirical evaluation of [10] |
| *UsersRegistry* [19] | Uses non-trivial API authentication |
| *PetStore* [9] | Used in the empirical evaluation of [24] |
| *SockShop* [38] | Used in the empirical evaluation of [1, 6] |
| *Openverse* [39] | Industry app, part of WordPress project |

## 4   Evaluation

In this section, we conduct a comprehensive evaluation of our methodology and its practical application in COTS. This assessment is twofold:

1. **Qualitative Analysis:** We aim to determine if the test-models defined by COPENAPI and the test drivers generated by COTS are effective in identifying logic-based faults.
2. **Quantitative Analysis:** We assess the effectiveness of our methodology in terms of code coverage, which is a standard metric for evaluating testing tools.

Given the novelty of our approach in user-written, model-based testing for REST-APIs, we establish our baseline comparison with: *i)* fully-automated REST testing approaches; and *ii)* manually crafted REST API tests.

### 4.1   Experiment Setup

To conduct such evaluation, and obtain the results presented in Sect. 4.2, we follow the preparatory steps illustrated in the rest of this section: we first select the artefacts, build their COPENAPI model, and we determine an adequate number of test runs per application.

*Artefact selection.* To conduct such an evaluation, we consider a sample of third-party applications listed in Table 1, satisfying the following criteria:

1. they must be open source, to facilitate reproducibility of our results;
2. they must have OpenAPI specifications, needed by COTS;
3. they should be non-trivial and thus representative of real-world RESTful applications;
4. ideally, they should include manually-written tests of their REST APIs, should be amenable to code coverage measurements, and should have been already used for evaluation in previous literature on testing.

All applications in Table 1 are open source, provide OpenAPI specifications, and are non-trivial; moreover, all of them (except *UsersRegistry* and *Openverse*) have been used in previous testing literature. The first 5 applications also satisfy the rest of our "ideal" criteria (item 4 above): they include handcrafted test suites for their REST APIs, and their architecture (consisting of a single Java-based executable) allow us to easily collect and analyse code coverage information using standard tools. Two other applications (*LanguageTool* and *PetStore*) also have similar architecture (allowing us to analyse code coverage), but do not include a test suite for their REST APIs.

*Building the* COPENAPI *Models.* As we are not the authors of the SUTs in Table 1, we were required to infer the usage of the SUTs as intended by their developers—in particular, what sequences of operations (request/responses) are valid, and how some requests may depend on others. As mentioned earlier, this information is typically only given informally in the application documentation, it may not be up-to-date, and is often omitted. Therefore, we often inferred this information by examining the pre-existing handcrafted tests, or by studying existing REST API clients. We also enriched our test models with application usage scenarios that, although possible, might have been overlooked in the existing handwritten tests. This way, our models may leverage SUT functionalities beyond those covered by the existing tests, and potentially reveal new faults in the SUTs.

*Example 3. FeaturesService* was the only application with documentation about some of its operations. For instance, the documentation says:

> *"The application should allow one to add a constraint such that when a feature requires another feature to be active, the latter feature cannot be deactivated without first deactivating the former."*

To test whether *FeaturesService* respects such a description, we formalise the COPENAPI model shown in Fig. 5 (abridged). This model specifies that the test driver should first add two features by invoking `addFeature`; their names are created by random generators and bound in variables *feat1* and *feat2*. Then, the test driver should invoke `addConstraint` specifying that *feat2* requires *feat1*. Finally, the test driver should attempt to delete *feat1* via `delFeature`, expecting the SUT to answer with a `400` "bad request" response. If the SUT was to answer with anything other than a `400`, (such as `200`, indicating that the operation was successful) it would indicate a fault in the logic of the SUT.          □

*Determining an Adequate Number of Test Runs.* After writing the COPENAPI model of each SUT and generating its test driver using COTS, we established an adequate number of test runs for each SUT. The optimal number of test runs is the smallest number that maximises *(1)* the number of discovered faults, *(2)* the level of code coverage, and *(3)* the level of model coverage. This number depends on the complexity of the SUT and the COPENAPI model in use. To determine the optimal number, we adopted an incremental approach: gradually

```
1  S_featuresService = rec X.(
2      +{ !addFeature(feat1: String(genFeatName)).?C201().
3              !addFeature(feat2: String(genFeatName)).?C201().
4                  !addConstraint(feat1, feat2).?C201().
5                      !delFeature(feat1).?C400().X,
6          ...
7        } )
```

**Fig. 5.** COPENAPI model for testing *FeaturesServices* (abridged), discussed in Example 3.

increase the number of test runs until the rate of bug discovery and code coverage plateau. This was determined by analysing the logs generated by COTS for each application and verifying that the model was being fully traversed.

### 4.2 Results

We now present our results, in terms of discovered faults (Sect. 4.2) and code coverage (Sect. 4.2).

**Discovered Faults.** Table 2 reports the number of faults discovered by our test models and the COTS-generated test drivers, categorised by the test oracle that detected the fault. We found 25 faults across the 9 selected applications and reported some of the most representative faults to the application maintainers; when we received a reply, in almost all cases the developers acknowledged that the faults were real bugs.[3]

Table 2 classifies the faults detected by COTS into logic-based and systematic categories. Logic-based faults, embedded in the SUT's logic, are discerned through testing specific operational sequences, while systematic faults, inherent in the SUT's overall function, can be identified by testing the requests in isolation of other requests. COTS successfully uncovered 10 logic-based and 15 systematic faults. Detecting the logic-based faults is attributed to COPENAPI models, which incorporate the SUT's domain knowledge, enabling COTS to navigate complex, interdependent request-response sequences. Such intricate sequences are unlikely to be deduced by fully-automated tools (e.g., Example 3), that, as mentioned earlier focus more on systematic faults. Furthermore, as Table 2 demonstrates, COTS effectively detects both fault types.

Some examples of logic-based faults follow:

– In the *PetClinic* model, we included an assertion to verify that a newly created resource can be subsequently retrieved from the SUT. However, COTS

---

[3] For *LanguageTool*, we reported 3 cases of *"500 Internal Server Error"* uncovered by our test model; such errors are considered faults in REST testing literature, but *LanguageTool* produces such errors to signal an unsupported functionality.

**Table 2.** Faults discovered by COTS test drivers and oracles, using our COpenAPI models.

ïž£

| Application | Bad status code | Bad response body | Assertion fail | Total | Logical errors | Systematic errors |
|---|---|---|---|---|---|---|
| *FeaturesService* | 6 | 0 | 0 | 6 (1, 0) | 2 | 4 |
| *RESTCountries* | 0 | 2 | 0 | 2 (2, 0) | 2 | 0 |
| *GestaoHospital* | 1 | 0 | 0 | 1 (1, 0) | 0 | 1 |
| *LanguageTool* | 3 | 0 | 0 | 3 (3, 3) | 0 | 3 |
| *PetClinic* | 5 | 0 | 2 | 7 (3, 3) | 4 | 3 |
| *UsersRegistry* | 1 | 0 | 0 | 1 (1, 1) | 0 | 1 |
| *PetStore* | 1 | 0 | 0 | 1 (1, 0) | 0 | 1 |
| *SockShop* | 1 | 1 | 0 | 2 (2, 0) | 1 | 1 |
| *Openverse* | 1 | 1 | 0 | 2 (2, 2) | 1 | 1 |
| **Total** | 19 | 4 | 2 | 25 (16, 6) | 10 | 15 |

discovered a fault: if a resource is updated, subsequent retrieval attempts fail, indicating an issue in the SUT's retrieval operation.

– Another fault in the *PetClinic* application involved inconsistent identification numbers. The SUT issues an ID number when a resource is created, but the list of resources returned shows different IDs than those originally assigned.
– The *FeaturesService* application did not adhere to the model shown in Example 3. It incorrectly permitted the deactivation of *feat1*, which should have been disallowed, and erroneously returned a successful HTTP status code of 204 instead of indicating failure.
– In the *SockShop* application, a significant fault was identified: the application fails to retrieve resources immediately after their creation. This fault was detected by first creating the resource and then attempting to retrieve it via the API sequence modelled via COpenAPI.

**Code Coverage.** Table 3 reports the size of the case studies,[4] and compares COTS model-based testing coverage results against handcrafted REST API test suites (when available as part of the selected case studies), and against fully-automated REST API testing.

The table shows the sizes (in lines of code) of our COpenAPI models against the handwritten REST API test suites. The COTS model sizes include the preambles, whereas the handwritten test sizes exclude their comments. Table 3

---

[4] The *LanguageTool* application has 5 API endpoints but only 2 are testable.

**Table 3.** Size information about the case studies, and comparison between COTS-based testing results (with COpenAPI models developed by us, executed for up to 1 min with default initial randomness seed), handwritten tests (part of the SUT source code), and Morest (average across 5 repetitions, each requiring 8 h of execution). The coverage is measured using JaCoCo 0.8.8 (https://www.eclemma.org/jacoco) with standard configuration; JaCoCo is also used in the Morest evaluation [24].

| Application | Application LOC | #endpoints | COpenAPI model size (LOC) | Manual tests size (LOC) | Number of COTS test runs | COTS tests line coverage | Manual tests line coverage | Morest avg. line coverage [24] |
|---|---|---|---|---|---|---|---|---|
| RESTCountries | 2409 | 22 | 247 | 300 | 100 | **1722** | 896 | |
| GestaoHospital | 4427 | 20 | 138 | 463 | 50 | **2857** | 2532 | |
| PetClinic | 10,416 | 35 | 225 | 1321 | 100 | 3099 | **3127** | |
| UsersRegistry | 5452 | 11 | 68 | 246 | 50 | **2035** | 1906 | |
| FeaturesService | 2026 | 18 | 98 | 377 | 150 | **1626** | 1576 | 360 |
| LanguageTool | 18,053 | 2 | 75 | | 20 | **4999** | | 935 |
| PetStore | 3693 | 20 | 111 | | 100 | **1987** | | 763.20 |
| SockShop | 3392 | 15 | 99 | | 50 | | | |
| OpenVerse | 7117 | 16 | 124 | | 1 | | | |

also shows the number of test runs performed using COTS and the number of lines covered by 1. COTS; 2. the applications' handwritten REST API tests; and 3. the fully-automated tool Morest [24,25]. We focus on Morest because a recent study [24] shows Morest to be superior to the other fully-automatic REST API testing tools in literature. Therefore, we selected Morest as representing the state-of-the-art, and included some of its case studies (*FeaturesService*, *LanguageTool* and *Petstore*) in our experiments.

In Table 3 we can observe that COpenAPI models are smaller than the handwritten tests provided by the SUTs—and yet, their code coverage is comparable, and higher for COTS in most cases. In the case of *RestCountries*, the COpenAPI model size is quite close to the handwritten tests size, but the COTS coverage is significantly higher. This suggests that a concise COpenAPI model can replace (part of) a larger handwritten test suite, allowing testers to "*concentrate on a (data) model and generation infrastructure instead of hand-crafting individual test*" [11].

With respect to the fully-automatic testing tool Morest, the benefits of developing a COpenAPI model are evident in the significantly higher line coverage achieved by COTS. In the three common applications we examined, COTS achieved 3 times as much coverage in *PetStore*, 4 times as much in *FeaturesService*, and 5 times as much in *LanguageTool*. Our results were obtained in a few seconds (once the COpenAPI model was developed), whereas the Morest coverage is averaged over 5 executions taking 8 h each. This suggests that, after the

initial investment of time needed to write a COpenAPI model, our model-based approach pays off in terms of coverage achieved and time saved over repeated executions.

### 4.3    Threats to Validity and Limitations

The time taken to develop COpenAPI models for these case studies varied from one application to another, depending on the complexity of the requests of the applications. On average, we estimate that it took us 30 h to complete each model—including the time we used to infer the intended usage of the application in question. We believe that a developer or tester more familiar with the application domain and requirements would write equivalent (or better) test models in less time.

*Internal Threats to Validity.* The method for determining optimal test runs could lead to imprecisions due to local maxima. Our lack of precise knowledge about expected request/response sequences necessitated designing COpenAPI models based on application tests and source code, potentially introducing bias (Table 3). Additionally, the comparison does not account for the qualitative distinction between creating concrete test cases and developing more abstract COpenAPI models, with the latter possibly being viewed as more challenging.

*External Threats to Validity.* Our case studies, though diverse, may not universally represent all application types, limiting generalizability (Sect. 4.1). We aligned our study with REST API testing discourse, selecting artefacts common in literature to minimize bias and ensure relevance. Despite these efforts, the specific selection could limit applicability, though it was necessary for comparability with existing research.

## 5    Related Work

In contrast to existing methodologies, our COpenAPI DSL allows for a more nuanced and comprehensive modelling of REST API behaviour. It enables the description of complex sequences and dependencies that are essential for thoroughly testing RESTful services. This approach is a significant departure from traditional models, which typically address simpler scenarios or focus on individual API operations in isolation.

We classify studies on testing of REST APIs into two broad categories: *model-based* and *fully automatic*.

**Model-Based Testing of REST APIs.** Our work enriches model-based testing of REST APIs by introducing a sophisticated DSL, COpenAPI, for detailed modeling and testing, diverging from existing approaches by supporting complex test sequences and state dependencies within a single model (§3). Unlike Chakrabarti *et al.* [7] and Fertig *et al.* [15], who focus on isolated sequences or operations, our DSL enables multiple test paths and value assertions, surpassing

the limitations of singular test sequence models and isolated API testing. Furthermore, while Seijas *et al.* [32] and Pinheiro *et al.* [30] contribute valuable perspectives on property-based and state machine-based modeling, they lack in addressing dynamic data generation and comprehensive evaluation. Francisco *et al.*'s [17] constraints-driven approach and Aichernig *et al.*'s [2] business rule modeling offer insights into precondition and postcondition specification, yet neither adequately tackles the RESTful API stateful interaction complexity our COpenAPI DSL is designed for. In contrast, TorXakis, another model-based testing tool, applied in [34] to the Dropbox synchronization protocol, does not support RESTful APIs and extending it would entail significant development, including the creation of an adapter for translating REST API requests and responses and aligning OpenAPI data types with those of TorXakis [34]. Our approach uniquely facilitates capturing intricate dependencies across REST API requests and responses, a critical aspect for thorough testing of sophisticated web services.

**Fully Automatic Testing of REST APIs.** This category includes a variety of testing approaches [13,31,33,43] and tools such as RestTestGen [37], RESTler [5], EvoMaster [4], QuickREST [21], RestCT [40] and Morest [24,25]. Starting from an API specification like OpenAPI [28], these tools automatically generate tests primarily to check the correctness of individual API responses. They mostly rely on random data generators or generators automatically derived from the specifications, which might not fully capture the intricacies of the API's behaviour. In contrast, in our work, the sequencing of valid requests and responses is explicitly specified in the model, enabling complex interactions with the SUT. This is complemented by our use of non-random, tester-assigned data generators in the COpenAPI model, which allows for a more targeted and effective testing approach. This method provides a stark contrast to AGORA's invariant detection-based approach [3] and the approach in [26] focussing on the specification of inter-parameter dependencies. While these methodologies offer robust solutions in their respective areas, they do not offer the same level of specificity in interaction sequencing and tailored data generation as our COpenAPI model-driven approach does.

## 6   Conclusion

We introduced COTS, a tool implementing a model-based testing approach for RESTful applications using COpenAPI DSL and OpenAPI specifications, which generates executable test drivers for model-based testing of web services. Our evaluation against third-party open source applications and comparison with existing automated REST API testing tools and manual test suites demonstrates the efficacy of COTS. It achieved comparable or better code coverage and fault detection with smaller, more manageable models, highlighting the potential for reduced effort in test creation and maintenance. Future work includes extending COTS to support other API specifications like GraphQL and gRPC, and enhancing COpenAPI's expressiveness by incorporating additional test oracles and timing constraints.

# References

1. Aderaldo, C.M., Mendonça, N.C., Pahl, C., Jamshidi, P.: Benchmark requirements for microservices architecture research. In: ECASE@ICSE, pp. 8–13. IEEE (2017)
2. Aichernig, B.K., Schumi, R.: Property-based testing of web services by deriving properties from business-rule models. Softw. Syst. Model. **18**(2), 889–911 (2019)
3. Alonso, J.C., Segura, S., Ruiz-Cortés, A.: AGORA: automated generation of test oracles for REST APIs. In: ISSTA, pp. 1018–1030. ACM (2023)
4. Arcuri, A.: Restful API automated test case generation. In: QRS, pp. 9–20. IEEE (2017). https://doi.org/10.1109/QRS.2017.11
5. Atlidakis, V., Godefroid, P., Polishchuk, M.: RESTler: stateful REST API fuzzing. In: ICSE, pp. 748–758. IEEE/ACM (2019)
6. Avritzer, A., Ferme, V., Janes, A., Russo, B., Schulz, H., van Hoorn, A.: A quantitative approach for the assessment of microservice architecture deployment alternatives by automated performance testing. In: Cuesta, C.E., Garlan, D., Pérez, J. (eds.) ECSA 2018. LNCS, vol. 11048, pp. 159–174. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-00761-4_11
7. Chakrabarti, S.K., Kumar, P.: Test-the-rest: an approach to testing restful webservices. In: 2009 Computation World: Future Computing, Service Computation, Cognitive, Adaptive, Content, Patterns, pp. 302–308 (2009). https://doi.org/10.1109/ComputationWorld.2009.116
8. Community, S.P.: Pet clinic application (2022). https://github.com/spring-petclinic/spring-petclinic-rest
9. Community, S.A.: Pet store application (2023). https://github.com/swagger-api/swagger-petstore
10. Corradini, D., Zampieri, A., Pasqua, M., Ceccato, M.: Empirical comparison of black-box test case generation tools for restful APIs. CoRR abs/2108.08196 (2021)
11. Dalal, S.R., et al.: Model-based testing in practice. In: ICSE, pp. 285–294. ACM (1999)
12. Dezani-Ciancaglini, M., de'Liguoro, U.: Sessions and session types: an overview. In: Laneve, C., Su, J. (eds.) WS-FM 2009. LNCS, vol. 6194, pp. 1–28. Springer, Heidelberg (2010). https://doi.org/10.1007/978-3-642-14458-5_1
13. Ed-Douibi, H., Izquierdo, J.L.C., Cabot, J.: Automatic generation of test cases for REST APIs: a specification-based approach. In: EDOC, pp. 181–190. IEEE Computer Society (2018)
14. noz Ferrara, J.M.: Features model microservice (2016). https://github.com/JavierMF/features-service

15. Fertig, T., Braun, P.: Model-driven testing of restful APIs. In: WWW (Companion Volume), pp. 1497–1502. ACM (2015)
16. Fielding, R.T., Taylor, R.N.: Architectural styles and the design of network-based software architectures. Ph.D. thesis, University of California, Irvine (2000)
17. Francisco, M.A., López, M., Ferreiro, H., Castro, L.M.: Turning web services descriptions into quickcheck models for automatic testing. In: Erlang Workshop, pp. 79–86. ACM (2013)
18. GraphQL (2023). https://graphql.org
19. Henrique, R.: Users registry application (2022). https://github.com/Throyer/springboot-api-rest-example
20. Hüttel, H., et al.: Foundations of session types and behavioural contracts. ACM Comput. Surv. **49**(1), 3:1–3:36 (2016). https://doi.org/10.1145/2873052
21. Karlsson, S., Causevic, A., Sundmark, D.: QuickREST: property-based test generation of OpenAPI-described restful APIs. In: ICST, pp. 131–141. IEEE (2020)
22. Kim, M., Xin, Q., Sinha, S., Orso, A.: Automated test generation for rest APIs: no time to rest yet. In: Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis, pp. 289–301 (2022)
23. Lenarduzzi, V., Daly, J., Martini, A., Panichella, S., Tamburri, D.A.: Toward a technical debt conceptualization for serverless computing. IEEE Softw. **38**(1), 40–47 (2021)
24. Liu, Y., et al.: Morest: model-based RESTful API testing with execution feedback. In: ICSE, pp. 1406–1417. ACM (2022)
25. Liu, Y., Li, Y., Liu, Y., Wan, R., Wu, R., Liu, Q.: Morest: industry practice of automatic restful API testing. In: ASE, pp. 138:1–138:5. ACM (2022)
26. Martin-Lopez, A., Segura, S., Müller, C., Ruiz-Cortés, A.: Specification and automated analysis of inter-parameter dependencies in web APIs. IEEE Trans. Serv. Comput. **15**(4), 2342–2355 (2022)
27. Matos, A.: REST Countries. https://restcountries.com
28. OpenAPI (2023). https://www.openapis.org
29. Language Tool Organisation: Language tool. https://github.com/languagetool-org/languagetool
30. Pinheiro, P., Endo, A.T., da Silva Simão, A.: Model-based testing of RESTful web services using UML protocol state machines (2013)
31. Segura, S., Parejo, J.A., Troya, J., Cortés, A.R.: Metamorphic testing of restful web APIs. IEEE Trans. Softw. Eng. **44**(11), 1083–1099 (2018)
32. Seijas, P.L., Li, H., Thompson, S.J.: Towards property-based testing of restful web services. In: Erlang Workshop, pp. 77–78. ACM (2013)
33. Stallenberg, D.M., Olsthoorn, M., Panichella, A.: Improving test case generation for REST APIs through hierarchical clustering. In: ASE, pp. 117–128. IEEE (2021)
34. Tretmans, G.J., Laar, M.: Model-based testing with TorXakis: the mysteries of dropbox revisited (2019)
35. de Vargas, V.P.: Gestao hospital system. https://github.com/ValchanOficial/GestaoHospital
36. Vasconcelos, V.T.: Fundamentals of session types. Inf. Comput. **217**, 52–70 (2012). https://doi.org/10.1016/j.ic.2012.05.002
37. Viglianisi, E., Dallago, M., Ceccato, M.: RESTTESTGEN: automated black-box testing of restful APIs. In: ICST, pp. 142–152. IEEE (2020)
38. Weaveworks: Sock shop (2021). https://github.com/microservices-demo/microservices-demo
39. WordPress Foundation: Openverse. https://wordpress.org/openverse

40. Wu, H., Xu, L., Niu, X., Nie, C.: Combinatorial testing of restful APIs. In: ICSE, pp. 426–437. ACM (2022)
41. Zhang, M., Arcuri, A.: Open problems in fuzzing RESTful APIs: a comparison of tools. CoRR abs/2205.05325v2 (2022). https://arxiv.org/abs/2205.05325v2. Version 2
42. Zhang, M., Marculescu, B., Arcuri, A.: Resource-based test case generation for restful web services. In: GECCO, pp. 1426–1434. ACM (2019)
43. Zhang, M., Marculescu, B., Arcuri, A.: Resource and dependency based test case generation for restful web services. Empir. Softw. Eng. **26**(4), 76 (2021)