

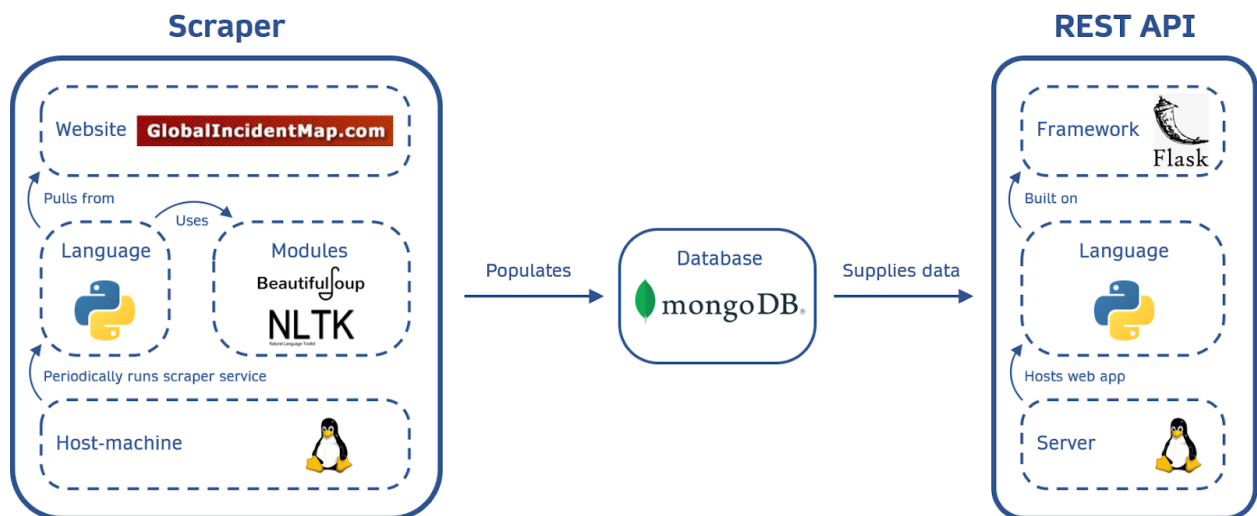
# Design Detail - Deliverable 1

## 1 Overall Design

### 1.1 Web Scraper

For the web scraper we will be using Python due to its ease of use. Being an interpreted language, developing a webscraper should be quick and simple. The primary module we will use is BeautifulSoup which allows us to parse and search html pages with ease. We will use other libraries and modules like regex, word2number and Natural Language Toolkit to extract more detailed information.

The site we are scraping ([globalincidentmap](http://globalincidentmap.com)) offers detailed reports sorted in chronological order. Consequently the scraper will be run daily parsing only data it has not seen before. It may be the case that there is more information to be found in the articles it links to. If this is the case then we will have the scraper explore those articles as well searching for keywords.



### 1.2 REST API

The API module will be developed using Python with Flask, a web app development framework. Using this framework a simple web application can easily be created that allows users to submit data to an endpoint of the website using a HTTP request. Each request will include a parameter that specifies what data the response should include. The web application will use the parameter to filter the records in the database and return the appropriate results.

The web scraper explained in 1.1, will populate a MongoDB database that the API will access for the data it will send as a response. MongoDB will be used since it stores data in a JSON format, the same format used throughout the application, making it easier to interact with. This web application will need to be hosted on a server, such as on a CSE machine, to provide the ability to run in web service mode.

## 1.3 Swagger

To aid with the use of our API by other teams we have decided to use Swagger. Swagger allows for clear documentation of our API, in an easy to read form. It also allows for someone to create dummy calls to our API through the Try It Out feature. Thus, another person is able to see how the calls work and what form they should be in. This aids with understanding as they are easily able to see what the calls should look like, thus easy integrate our API into their web app.

## 1.4 Web App

### 1.4.1 Frontend

The web app is where we can showcase our API data, so it is important to choose a frontend tech stack that enables us to quickly create a working app, while at the same time allowing for complex component design that can show off our features.

The framework we chose to use was React. This is for two reasons. Firstly, many of us have industry experience working with React, and so we can write the website at a faster pace. React also allows for higher complexity with its component system, enabling our team to go above and beyond, without being hindered by the framework.

To communicate with the backend, we are using Apollo GraphQL. This is widely accepted to be the most straightforward way of communicating with the backend through a single endpoint, and allows our users to offload the processing of data to the database. This means that our website uses less power and loads at a faster rate than comparable websites using standard fetch requests. It also makes our endpoint much easier to use, as they only have to understand the schema for a single endpoint.

There are three supporting libraries that we are utilising: Typescript, ESLint and Docker.

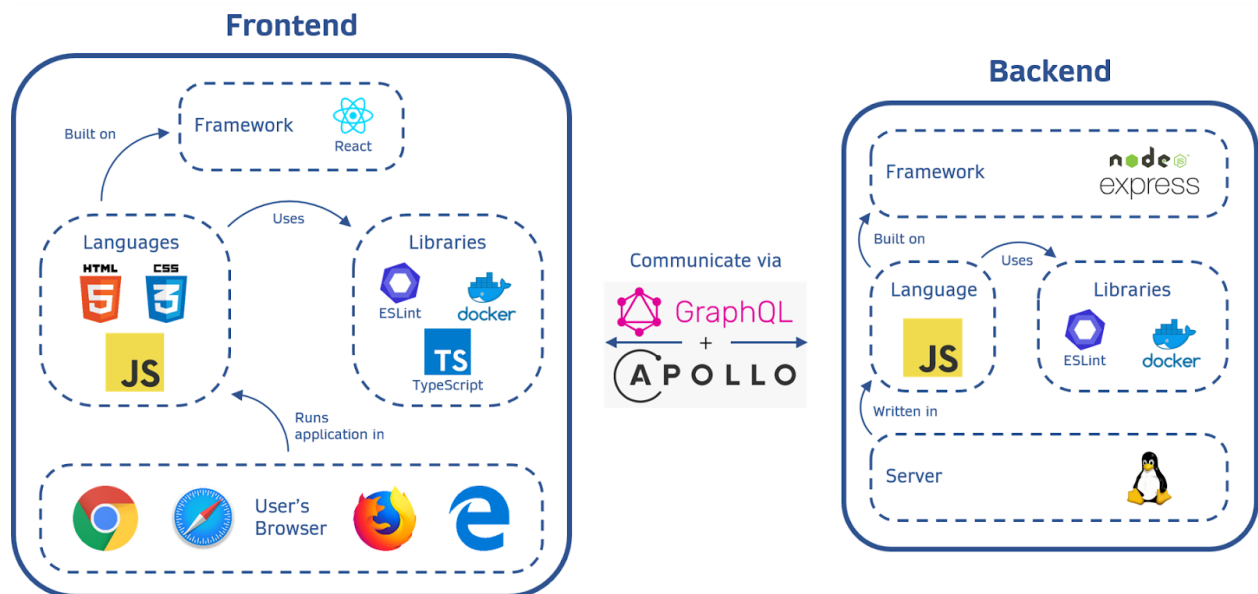
Typescript allows us to typecast our javascript code, essentially reducing the time spent tracking bugs caused by javascript's dynamic typing system.

ESLint allows us to standardise our code style, reducing the time it takes for new developers to understand what is going on. It also does this automatically, so no extra time is spent on formatting.

Docker allows us to standardise the testing and development of the frontend. This means that everyone who runs our code should have identical outputs, allowing for easier bug fixing. It also automates a lot of the setup, allowing more devs to quickly become active in the development workflow when needed.

## 1.4.2 Backend

Our backend will be quite lightweight, as most complex analysis and data will be retrieved from our public API. Due to this, we decided to use NodeJS Express, which is known for its flexibility and ease of use. Express also has a GraphQL library which we can use to serve our services from a single endpoint. Just like the frontend, we will use ESLint and Docker. Docker in particular will allow us to quickly run a test backend and frontend instance in conjunction, improving our testing capabilities.



## 2 API Design

### 2.1 Endpoint Design

Clients can interact with the API through a single endpoint via HTTP requests. These requests will use parameters to contain all the necessary information required by the API to calculate and send back the appropriate results. Depending on the data contained in the parameter, the API will filter the data stored in the database and return the requested data as a JSON object. The parameter must adhere to a strict structure made up of 4 pieces of information represented in the form: {"start\_date", "end\_date", "key\_terms", "location"}, where start\_date and end\_date are in the form "yyyy-MM-ddTHH:mm:ss". If the parameter does not adhere to the required structure

then the API will return an error result with an appropriate error code. Below is a table describing the possible response codes and the reason for them.

Response Code	Reason
200	Everything is ok and correct result is returned
400	There is an issue with the parameter
404	No data matching the parameter requirements could be found
405	The request method used is not allowed
500	The API had an issue processing the request

## 2.2 Example Interactions

Request Method	Parameters	Response	Test
GET	?start_date=2015-01-01 T00:00:00 &end_date=2018-01-23 T19:37:12 &key_terms=swineflu &location=punjab	{"date": "2015-02-15 17:33:00", "country": "Australia", "url": "http://www.abc.net.au/news/2015-02-16/pig-brucellosis/6116970", "coord": "-33.864, 151.205", "city": "New South Wales", "disease": "Brucellosis", "description": "[ABC.net.au]\u00a0 AUSTRALIA :: Call for pig dogs to be tested for brucellosis\n\n\"Veterinarians in the Hunter Valley and North West regions of New South Wales are calling on pig dog owners to be on the lookout for signs of brucellosis after several dogs tested positive for the disease.\"\n\nRead Full Article At :: http://www.abc.net.au/news/2015-02-16/pig-brucellosis/6116970"}	Normal Request
GET	?end_date=2018-01-23 T19:37:12 &key_terms=swineflu &location=punjab	400 status code with no JSON object returned due to the parameter not meeting the required structure	Missing Start Date

GET	?start_date=2015-10-01 T08:45:10 &end_date=2016-11-01 T19:37:12 &key_terms=covid19 &location=sydney	404 status code with no JSON object returned due to the API not having any data that meets the parameter requirements	No Data
POST	?start_date=2015-01-01 T00:00:00 &end_date=2018-01-23 T19:37:12 &key_terms=swineflu &location=punjab	405 status code with no JSON object returned due to the request method not being GET	Post Request
GET	This could happen independent of the supplied parameters, as a 500 status code is a catch for uncaught errors in the server.	500 status code with no JSON object returned due to an API error, possibly due to a missed edge-case	N/A

## 3 Implementation Choices

### 3.1 Implementation Language

The main languages that we have decided to go with are Python and JavaScript.

#### 3.1.1 Python

We are using Python for the web scraper and RESTful API due to the simplicity of the language. It allows us to focus on the implementation of the services, rather than getting bogged down in the complex structures that other languages, such as C or an Object Oriented language such as Java, poses. This also helps to increase the reliability of our code, which allows for seamless collaboration throughout the team.

Also, Python allows us to use a framework such as Flask, which handles a large amount of the overhead for creating a Web API. This allows us to yet again focus on the implementation of the API, increasing our productivity.

We understand that the choice to go with Python will reduce the speed of our program, as Python is inherently slower than a language like C. However, the benefits of increased productivity and readability outway the negative performance effects.

### 3.1.2 JavaScript

We have decided to use JavaScript, coupled with HTML and CSS, as our main frontend language as it has become an industry standard for web applications. Also, we are able to use frameworks such as React to help with the development of our frontend. React offers standard conventions and structuring that will help in the continued development of our web application. It will allow us to add new features and pages with ease, in a way that appears integrated with the rest of the site.

We will also be using NodeJS as our backend for the web application, coupled with Express, as this is a well used tech stack. The simplicity of Express will allow for the backend to be created simply and efficiently, and will also help in adding extra features as the project goes on.

## 3.2 Development and Deployment Environment

For our deployment environment we have decided to go with Linux. The reason we are using Linux is that it is lightweight and simple to set up for the small services that we are creating. Also, it is a common choice for services like ours so there is a large amount of documentation and help online, so we will be able to easily setup the services and diagnose issues when they arise.

However, all our frameworks work cross platform so during development we can all use them independent of our own environments. This speeds up development, as we do not have to individually configure our environments to help run this process. And with the introduction of Docker we will be able to perform the setup for our environments quickly.

We have decided to deploy our project onto Google Cloud Platform. The reason is that its highly versatile in what can run on the service, as it has both virtual machines for running the web scraper and app deployment for our API and Web App. It also supports linux deployment, which is our chosen environment.

## 3.3 Libraries

### 3.3.1 Typescript

Typescript is an open source library that adds typing to traditional javascript. This improves clarity, as a developer can discern more information from any given variable, and also reduces room for error when variables are misinterpreted. The overarching reason we added typescript was to improve our quality, and pace, of code written.

### 3.3.2 Apollo

Apollo is the industry-standard GraphQL implementation for services written in javascript. Instead of using a traditional REST service, we can access the backend through a single endpoint, allowing our users to offload most computation to our services. This means that our app uses less power and loads at a faster rate than comparable websites using standard fetch requests. It also makes our endpoint much easier to use, as they only have to understand the schema for a single endpoint.

### 3.3.3 Beautiful Soup

Beautiful Soup is an open source library which allows aids in the extraction of data from HTML files. It creates a parse tree of the file which can be easily navigated, searched and modified in idiomatic ways. This is incredibly useful for scraping data off websites as it grants the ability to use extract data based on its HTML structure.

### 3.3.4 Natural Language Toolkit (NLTK)

NLTK is a collection of libraries and programs which help programs reason about natural language texts. These include tagging, machine learning, parsing, probability, tokenization, stemming and chunking of texts. NLTK will be used to find more detailed information from the articles and descriptions.

### 3.3.5 Requests

An open source library that makes sending HTTP requests very easy, abstracting most of the required functionality required with useful methods. Using this library's 'get' method along with the API's URL and a JSON object parameter, a HTTP GET request will be sent to the API and a data object returned that can be converted into a JSON object using the library's 'json' method, all in one convenient line.

### 3.3.6 Flask-PyMongo

Flask-PyMongo is an open source library that bridges Flask and PyMongo and provides some convenience helpers. PyMongo provides the connection to the MongoDB database used by the API.

# Design Details - Deliverable 2

## 1 Final Architecture

The final architecture that we developed is very similar to the architecture that we planned out in deliverable 1. However, there were some changes during implementation, such as additional features to improve the user experience and make it easier for us to monitor the API. The details of each component's architecture are described below.

We used MongoDB as our database, forming the bridge between the web scraper and the API. This allowed for data to be kept in a useful format, json. The web scraper could format the data it gets from the website into a JSON object and easily pass it into the database. From there the API could retrieve the data and perform little processing on it before sending it in the HTTP response.

The web scraper's final architecture is very similar to what was described in deliverable 1. It uses beautiful soup to parse HTML and regex for keyword searches and string replacements. Finally pymongo is used to push suitable data to MongoDB from which the api will query.

The REST API uses Python and Flask to create our functional endpoint allowing users to make GET requests for our data. An additional module used was graphene, which allowed us to implement an optional graphql endpoint for any users who preferred non-REST api's. The format of queries and examples are documented in Swagger. Another addition was an endpoint that allows us to visually monitor all API requests and their performance. This was a huge improvement as it provides a graphical way to understand the success of our API rather than using the large amounts of text provided by log files.

Finally, all of the above services are hosted using the Google Cloud Platform, a single solution to meet all our web hosting needs. It provided capability to both host the API as an application, and a virtual machine through which we were able to run the web scraper remotely. This platform also provided a convenient location for the API logs, enabling us to search through the logs, finding each user's request to the API and any errors that occurred. An example of the logs can be found in the image below.



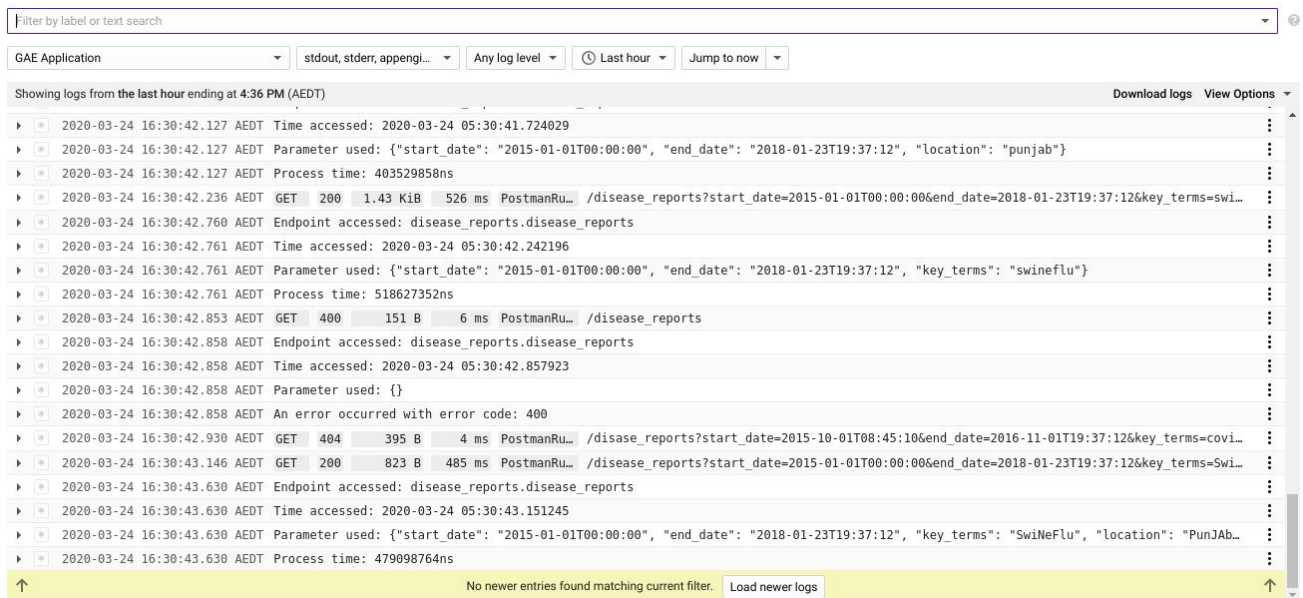
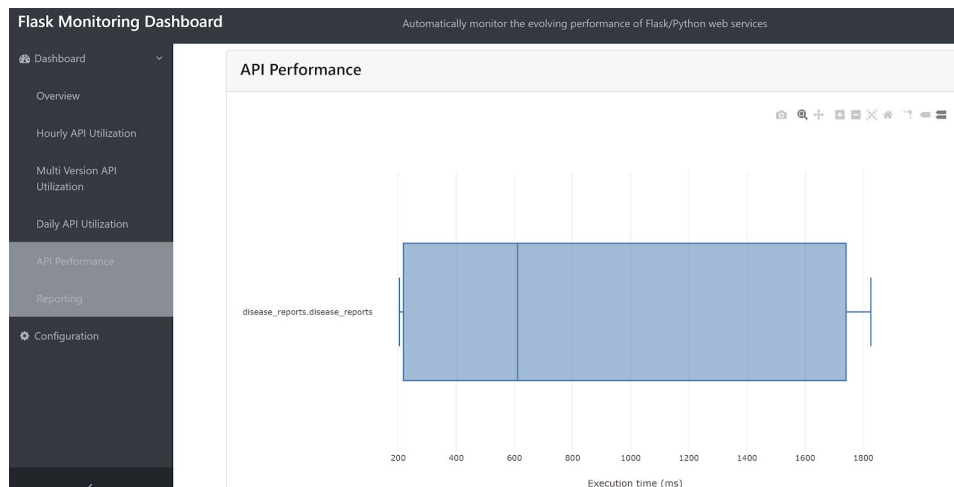


Image of Google Cloud Platforms logs

## 2 Choice of Implementation

### 2.1 REST API

As planned, the final API was developed using Python with the web app development framework Flask. This allowed for a quick development and deployment as the design was kept simple and small scale with one endpoint for functionality. What did change however was the addition of a new endpoint used to monitor resource utilisation visually, as shown in the image below. This won't affect how users interact with our API but will allow us to understand how they do and how well our API is performing.



Screenshot of Resource Utilisation Dashboard

The functional endpoint still involves users submitting parameters however they are no longer in the form of a JSON object and are instead passed as url query arguments. This made it easier to test as the API could be queried directly in the browser rather than using some testing service, an example can be seen in the image below. It also meant that our API adhered to industry standards since we are using the GET HTTP request and no longer sending a data object which would only be appropriate with a POST HTTP request.

## Example of Testing API Directly in the Browser

## 2.2 Web Scraper

## 2.3 Testing

## 3 Challenges and Shortcomings

### 3.1 Web Scraper

The site ([globalincidentmap](http://globalincidentmap.com)) we are parsing has tens of thousands of records which can be accessed incrementally by changing the url. Unfortunately some of these records are empty and consequently have to be ignored. This is trivial however it makes knowing when the scraper has reached the latest record difficult. To get around this, the scraper runs hourly picking up on the last successfully scraped record and stops after seeing 30 consecutive empty records or after scraping 500 records. This allows the scraper to be constantly up to date while being able move past many consecutive empty records.



Image of a blank/empty record

The main issue we faced writing the web scraper was converting the event type field of the aforementioned records into a disease or syndrome as described in the spec. Ebola and Murburg were both part of the same event type and consequently had to be seperated. Also all types of influenza were combined into a single category. These issues were for the most part surmountable by searching for specific keywords. Other event types had to be completely ignored such as “Ricin” and “Suspicious or Threatening Powder” due to their being irrelevant or off topic (There were many diseases which primarily exist in animals).

### 3.2 REST API

One of the most challenging aspects of implementing the API was figuring out how to query the database. This was particularly difficult due to the significant lack of documentation for PyMongo, the python library used to interact with our MongoDB database. For example, it was problematic finding how to query the database for records that are within a particular date range. To solve this issue, we had to make educated guesses based on the limited documentation and adjust based on the results until we eventually made it work.

Another problem we had to overcome was manipulating the parameters into a form that was suitable for querying the database. Thankfully, Python provides a wide range of libraries and internal functionality for manipulating strings, so with a little research the problem of string

manipulation was overcome. However, there remained the issue of understanding what format the data needed to be manipulated into, so that it matched the format of the database. This required communication with the team mate developing the web scraper as they made decisions of what the database format was.

As we had not had prior experience developing our own REST API's, the domain of possible cases was unknown to us. Through a process of trial and error we began to uncover the possible edge cases that could be presented. With use of our tests we were able to find and solve these edge cases effectively.

### 3.3 Swagger

The Swagger documentation of our API needed to clearly convey to users the structure and purpose of both the inputs and outputs of our system. Swagger is well suited to the latter, providing a schema section in which we can comprehensively define the structure of variables used, with examples, regex formats for strings, and lists of properties for objects. These schemas can then be referenced elsewhere, and Swagger automatically generates expandable examples for outputs with all the detail of the schema included.

However, when schemas are referenced from inputs, Swagger runs into two issues. The first is that the examples pulled from the schemas don't provide an immediately meaningful set of placeholders to the "Try it out" feature, which allows Swagger to make calls to our APIs, so we defined specific examples for each of the input values.

The second issue was that the schema details don't get shown with the inputs, making it difficult to tell what the structure of the inputs should be without reading through all the schemas to determine which applies. Anticipating that Swagger had some way of allowing us to reference this information, we went looking through the documentation for the Swagger application itself, and eventually found that the feature had been only recently programmed (just 2 weeks before the deliverable's deadline) and was scheduled to be included only in the next release. To ensure that users were still able to understand what each input should look like, we instead added descriptions and examples for the forms of each of the inputs to their descriptions, and where applicable directed users to the schema for the input's type for more information.

### 3.4 Testing

The original plan with the testing suite was to write a simple python script that would grab requests from text files and compare the response with an expected response. However, after complications with writing this small script, we realised that the response was more nuanced than hoped. A different order in the json object could cause the test to fail when it should pass.

The solution we came to here was to move all our tests to Postman. Postman has already handled all of these complications, and allows us to create more comprehensive tests much

easier. Also, we are able to run postman tests on both the command line and the postman client, giving other teams a choice with how they want to interpret our tests.

See the testing document

(<https://github.com/z5122506/SENG3011-calmClams/blob/master/Reports/Testing%20Documentation.pdf>) for testing details.

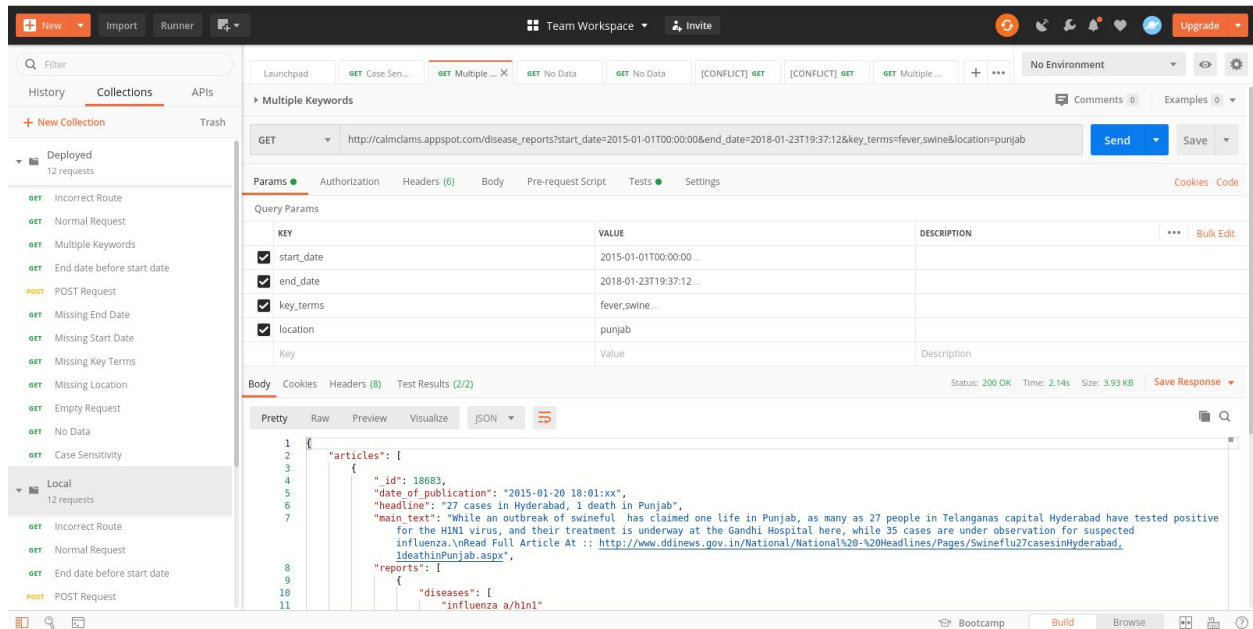


Image of our test suite in Postman