

# Artifact - Historia: Refuting Callback Reachability with Message-History Logics

September 9, 2023

## 1 Introduction

This document explains the scripts in the artifact to reproduce the results in the Historia paper [2]. In order to reproduce the subset listed later, an x86 Linux machine with 16GB of ram is needed (8GB in addition to the ram for the OS and other processes). The goal of this document is to first give a set of instructions for reproducing the experimental results and then give a technical explanation of how the implementation connects to the technical contributions.

When using HISTORIA on its own, the input is a compiled Android application in the form of an APK file, a CBCFTL specification defined in a JSON format, and an assertion at a location defined by a JSON data structure. Outputs are either “safe” or “alarm” in which case an abstract message history witnessing the alarm will be available.

A full tutorial on using HISTORIA may be found in the notebook `HistoriaExampleAndExplanation.ipynb` after running the Docker container. This document will focus on reproducing the results, but we encourage referencing of the descriptions in the tutorial notebook as needed.

## 2 Prerequisites - Running the Historia Docker Container

We have configured the experiments to be run within two Docker containers provided with this artifact. These Docker files may be found in the root directory of this archive and is labeled `historia.docker`. Please follow the instructions to install docker from <https://docs.docker.com/engine/install/>. Docker must be installed such that the daemon is running in root mode otherwise the notebooks shown in Figure 1 may not be visible (see <https://docs.docker.com/engine/install/linux-postinstall/>).

Importing the docker containers can be done with the following command. Please note that the exported docker container files `historia.docker` and `historia_postgres.docker` are contained in the root directory of the archive.

```
tar -xvf Historia.tgz
cd Historia
docker load < historia.docker
docker load < historia_postgres.docker
```

In order to start the Docker containers, run the following command in the root directory. This will start at least two docker containers.

```
docker compose -f docker-compose.yml up --scale worker=0
```

All subsequent steps may be done through the web interface at <http://localhost:8888>. Opening this URL should show a Jupyter notebook as shown by Figure 1.

When finished, the following commands will close down the containers and remove them.

```
docker compose down
docker compose rm
```

**System Requirements** We have split the instructions so that a subset of the experiments may be run on a reasonable laptop and included instructions to run the full experiments if access to a server is available. The instructions labeled “minimal resources” may be run on a reasonably modern laptop that has 8GB of ram in addition to ram used by other processes. If you are running on a system with more than 16 GB of ram, you may want to add additional workers by increasing the `--scale worker=` parameter. Each worker takes an additional 8GB of ram.

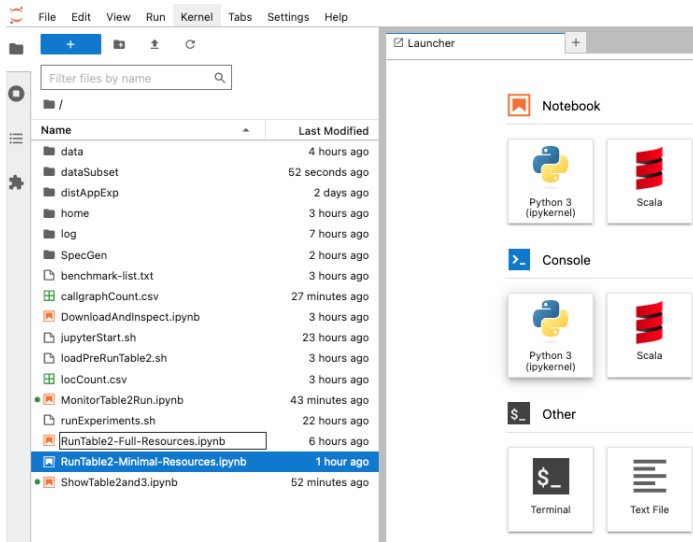


Figure 1: Your web browser should look like this if the Docker images started correctly.

Instructions labeled with "full resources" can be run on more powerful systems. We used multiple servers with an AMD EPYC 7763 64-Core Processor and 256GB of ram. Usage of these resources was granted to our project by Chameleon Cloud [1].

### 3 Reproducing the Historia Results

Each subsection here corresponds to a table in the evaluation section. We have labeled the subsections first with the research question ( **RQ1** or **RQ2** ), the table number (Table 1, Table 2, Table 3), and finally, the resource requirements as described earlier (minimal resources or full resources).

#### 3.1 RQ1 - Table 1 - minimal resources

This section shows how to reproduce table 1 except for row 4 fix which requires the full resources version. To run the experiments for the first table, open the terminal in the Jupyter notebook listed earlier and run the following command:

```
./runExperiments.sh
```

The experiments in Table 1 are run as a Scala unit tests defined in the file:

```
implementation/src/test/scala/edu/colorado/plv/bounder/symbolicexecutor/Experiments.scala.
```

Each benchmark is labeled by `test([description])` where [description] has a row number and short english description. The source code for each benchmark is a string stored in a variable named `src`. This source code is compiled into the APK automatically when running the unit test.

The set of CBCFTL specifications for each row are defined by row in the `ExperimentSpecs` object within `Experiments.scala`. These specifications may be found in:

```
src/main/scala/edu/colorado/plv/bounder/lifestate/Specification.scala.
```

The output will be printed to the screen in a log format. Each row of Table 1 may be found labeled by "Row" followed by the row number and version (i.e. "bug" or "fix"). For example, the buggy version of `getAct`<sup>[3]</sup> is "Row 1 bug" and the fixed version is "Row 1 fix".

Below is an explanation of each column in the table:

1. Pattern cb,ret - The number of callbacks and returns in the bug pattern: Integer labeled "cbSize".
2. Pattern ci - the number of callins in the bug pattern: Integer labeled "syntCi".
3. Historia specs - number of specs written for the benchmark. Integer labeled "spec count". The specific specs may be found in `Experiments.scala` and `Specification.scala` as described earlier.

4. Historia cb - number of callbacks that may be matched by the CBCFTL spec: Integer labeled “matchedCb”.
5. Historia cbret - number of callback returns that may be matched by the CBCFTL spec: Integer labeled “matchedCbRet”.
6. Historia ci - number of callins matched by the spec: Integer labeled “matchedSyntCi”.
7. Historia time - runtime:
8. Historia res - The result of the analysis: There will be several rows of text displaying the data from the table. The result of verification labeled with “actual:” and may say “Witnessed” (⊕), “Timeout” (⊗), or “Proven” (⊘). For reference, it also prints the expected result after “expected:”. If the actual and expected results differ, the unit test fails.

### 3.2 RQ2 - Table 2 and Table 3 - Viewing results

Since the full experiments take a long time to run, start by viewing our results in the Jupyter interface. During the experiments, jobs and results are stored in a database in the `historia_postgres` docker container started by `docker compose`. To load the pre-computed results, open a terminal in the Jupyter notebook and run the following script to load the final database:

```
./loadPreRunTable2.sh
```

The results can be viewed by opening the notebook `ShowTable2and3.ipynb` shown on the left. After opening, select “Kernel” -> “Restart Kernel and Run All Cells”. This notebook reads the results from the database and displays the Historia results for Table 2 and Table 3 (Shown under heading titles “Table 1” and “Table 2”). Additional statistics about the run may be viewed in this notebook as well. Statistics about the overall runtime may be found toward the end of the notebook `MonitorTable2Run.ipynb`.

### 3.3 RQ2 - Table 2 - minimal resources

The full set of 1090 locations can take an extremely long time to run (over 15 days of CPU time). Therefore, we recommend running a subset.

The notebook `RunTable2-Minimal-Resources.ipynb` will run by clearing the database and uploading the locations in the “dataSubset” directory. These locations will be processed by the Historia worker. The directory “dataSubset” contains 10 locations that don’t use too much ram and don’t time out. Locations are stored as “.json” files. Running should take around 20 minutes depending on the system.

Progress can be tracked through the `MonitorTable2Run.ipynb` notebook. Each time a cell in this notebook is run using “shift-enter”, the output will be updated. When the `completed_jobs` row under the “Track Total Jobs” heading reaches 10, the subset of locations is completed running.

Once completed, the last cell of the notebook `ShowTable2and3.ipynb` can compare the results just generated against our results. It will print out the path to the config file, the result from the most recent run, and finally, the result from the run used for the paper. Please note that memory errors often appear as the process simply being killed and generally do not show an error message.

**Changing the subset** The subset of apps may be changed to a random selection by deleting the contents of the “dataSubset” directory and re-running the last cell of `ShowTable2and3.ipynb`. Please note that many locations will not complete on the 8 GB of ram configured by the docker compose file (see the later instructions to remove this restriction).

**Call Graph Sizes** Sizes for the call graphs were computed by adding some code to count app methods after graph generation in both Flowdroid (line 1473 in `SetupApplication.java`) and HISTORIA. The call graph was then filtered for methods under the package declared by the application’s manifest. We recorded these call graph sizes in the file `callgraphCount.csv`. The version of Flowdroid we used for this comparison may be found in `/home/FlowDroid`.

### 3.4 RQ2 - Table 3 - minimal resources

The rows in Table 3 may be run using the `ShowTable2and3.ipynb` notebook. Under the heading “Table 3 Data - Write Specifications For Individual Samples” each cell (e.g. marked by “Specify\_0”) corresponds to one row in Table 3.

The output of the cell will contain the name of the app in the row. The configuration file used as input, and the directory where results are written. For example:

```
home/notebooks/SpecGen/com.darshancomputing.BatteryIndicatorPro/SensitiveDerefCallinCaused
/0/config.json
App Name: BatteryBot
out directory: /home/notebooks/SpecGen/com.darshancomputing.BatteryIndicatorPro/
SensitiveDerefCallinCaused/0
apk path: /fdroid/com.darshancomputing.BatteryIndicatorPro/12.0.0/apk/com.darshancomputing.
BatteryIndicatorPro\_26016.apk
```

Since we set the timeout to one hour, we recommend only running the rows that are shown to complete in the table. These rows can be run by uncommenting the call to `runHistoriaWithSpec()` at the bottom of the cell. For example, the `//` can be removed from the second two lines of BatteryBot:

```
//Uncomment the following two lines to run Historia on this sample from Table 3
// val result = Specify_1.runHistoriaWithSpec()
// println(s"Result Summary: ...")
```

Making it look like:

```
//Uncomment the following two lines to run Historia on this sample from Table 3
val result = Specify_1.runHistoriaWithSpec()
println(s"Result Summary: ...")
```

Afterwards, running the cell will take slightly longer than before and produce the “res” column of Table 3. For example, the output may look like:

```
Result Summary: Proven
```

### 3.5 RQ1 - Table 1 - full resources

In order to run the full version of Table 1, a system with at least 120G of ram must be used. To run, restart the Docker container using the following command:

```
docker compose -f docker-compose-full-resources.yml up
```

Then, repeat the instructions from the “minimal resources” section. This should take about an hour to run on a comparable system.

If this completes successfully, the output should be the same as earlier but will include “Row 4 fix”.

### 3.6 RQ2 - Table 2 - full resources

In order to run the full version of Table 2, we strongly suggest using multiple servers. You can use the `MonitorTable2Run.ipynb` notebook to estimate the expected remaining runtime after a few jobs have completed (See the “Estimate Time Until Completion of All Jobs” cell).

The Docker compose file `docker-compose-full-resources.yml` shows how to start and connect worker containers. Scaling to multiple systems can be done with Docker swarm <https://docs.docker.com/engine/swarm/> or simple SSH tunnels connecting port 5432 from the workers to the Postgres database container. Once the workers and server are configured, the full experiments may be run through the `RunTable2-Full-Resources.ipynb` notebook. To add more workers on a single system, the `--scale worker=n` flag may be added to the docker command.

The results may be viewed the same as was explained in the “minimal resources” section earlier.

### 3.7 RQ2 - Table 3 - full resources

The instructions for Table 3 may also be run under the full resources version of the docker containers, but the results are similar.

## References

- [1] Kate Keahey, Jason Anderson, Zhuo Zhen, Pierre Riteau, Paul Ruth, Dan Stanzione, Mert Cevik, Jacob Colleran, Haryadi S. Gunawi, Cody Hammock, Joe Mambretti, Alexander Barnes, François Halbach, Alex Rocha, and Joe Stubbs. 2020. Lessons Learned from the Chameleon Testbed. In *Proceedings of the 2020 USENIX Annual Technical Conference (USENIX ATC '20)*.
- [2] Shawn Meier, Sergio Mover, Gowtham Kaki, and Bor-Yuh Evan Chang. 2023. Historia: Refuting Callback Reachability with Message-History Logics.. In *Object-Oriented Programming Systems, Languages, and Applications (OOPSLA)*.