

# DisProTrack: Distributed Provenance Tracking over Serverless Applications

## Abstract

Provenance tracking, crucial for debugging system vulnerabilities, is widely employed. Existing models focus on monolithic applications, but modern DevOps-based service architectures pose challenges. This paper introduces a novel Universal Provenance Graph (UPG) approach, utilizing a Loadable Kernel Module (LKM) for runtime unit identification and a log optimization method for effective provenance tracking in serverless architectures. Evaluation with various benchmarked serverless applications demonstrates the model's effectiveness.

## Provenance Graphs

Provenance data is the metadata of a process capturing origin and modification details throughout its lifecycle. The resulting Provenance Graph is a causal graph depicting dependencies between system subjects (e.g., processes) and objects (e.g., files, network sockets).

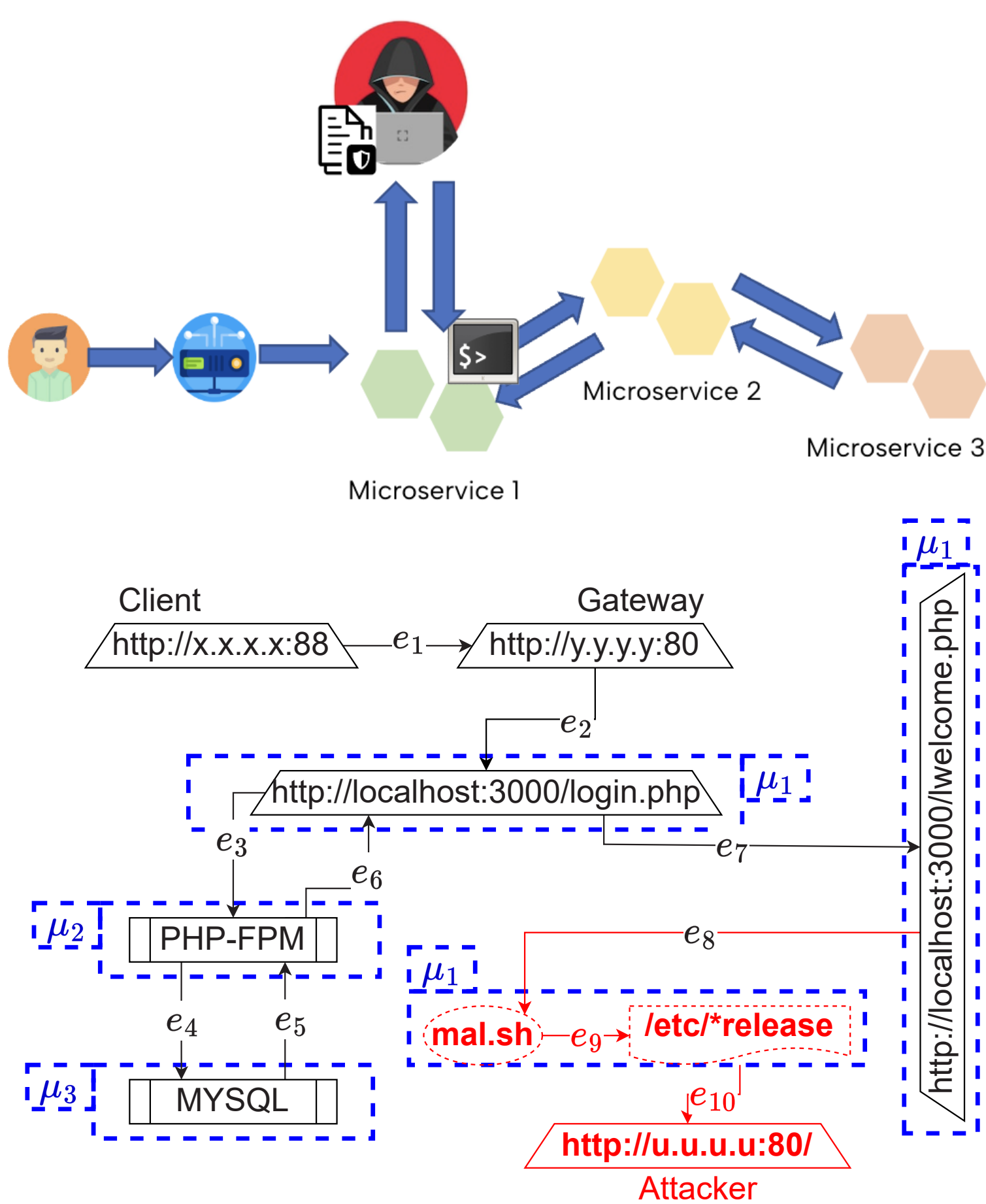


Figure 1. The Top Figure shows an Adversarial Model for Confidential Data Theft, and the bottom figure shows the UPG generated for the same.

### Log Message Generating Function

It is a library function used for printing a LMS in either terminals, specific log store files, or log databases. For example, consider the following C code snippet with a popular logging library Log4C. `log4c_category_log(NULL, LOG4C_PRIORITY_ERROR, "Hello World!");` Here the LMS is Hello World! and the LMGF is `log4c_category_log`.

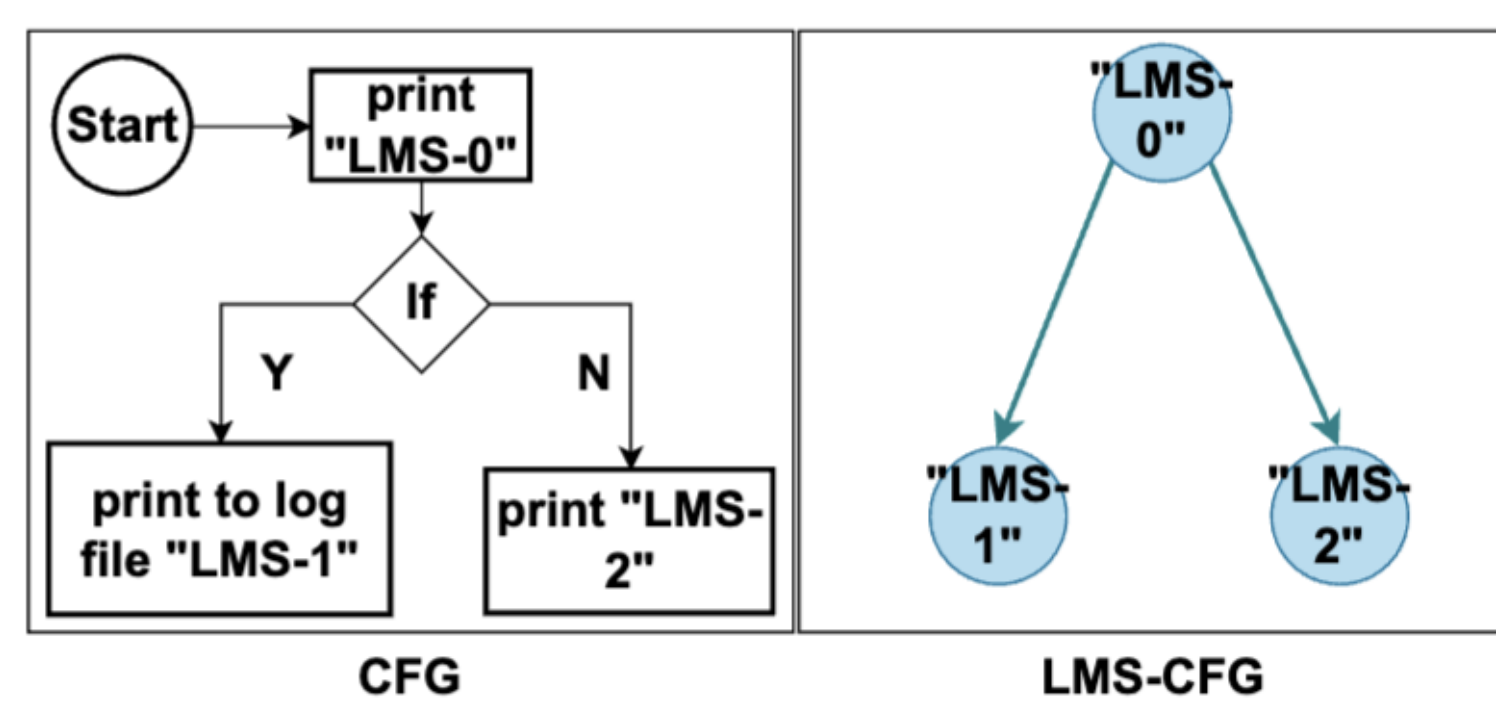


Figure 2. Example of CFG and LMS-CFG

### UPG Construction - Challenges

1. Combining application logs from different micro-services: Log messages formatting, timestamp formatting, process descriptors vary
2. Combining the system log with the application logs: Container-based sandboxing shares same pid namespace
3. Identification of execution units
4. Dependency explosion and handling confounding root causes: Reverse query results more than one root cause

## Contact Me

DisProTrack GitHub:<https://github.com/usatpath01/DisProTrack>  
My Homepage: <https://usatpath01.github.io/>

## DisProTrack Overview

### Contributions:

1. Design of the UPG from application and system logs: Static analyzer module generates the application-specific Log Message String-Control Flow Graph (LMS-CFG) from the application binaries which provides a profile of the application.
2. Runtime execution unit identification: Developed a Linux LKM that can intercept the system calls generated during execution time to identify the semantic relationship between the system logs and the application logs.
3. Utilization of Regular Expression to improve search efficacy: Instead of storing the raw log messages in the UPG, we propose conversion and storage of an equivalent regular expression.

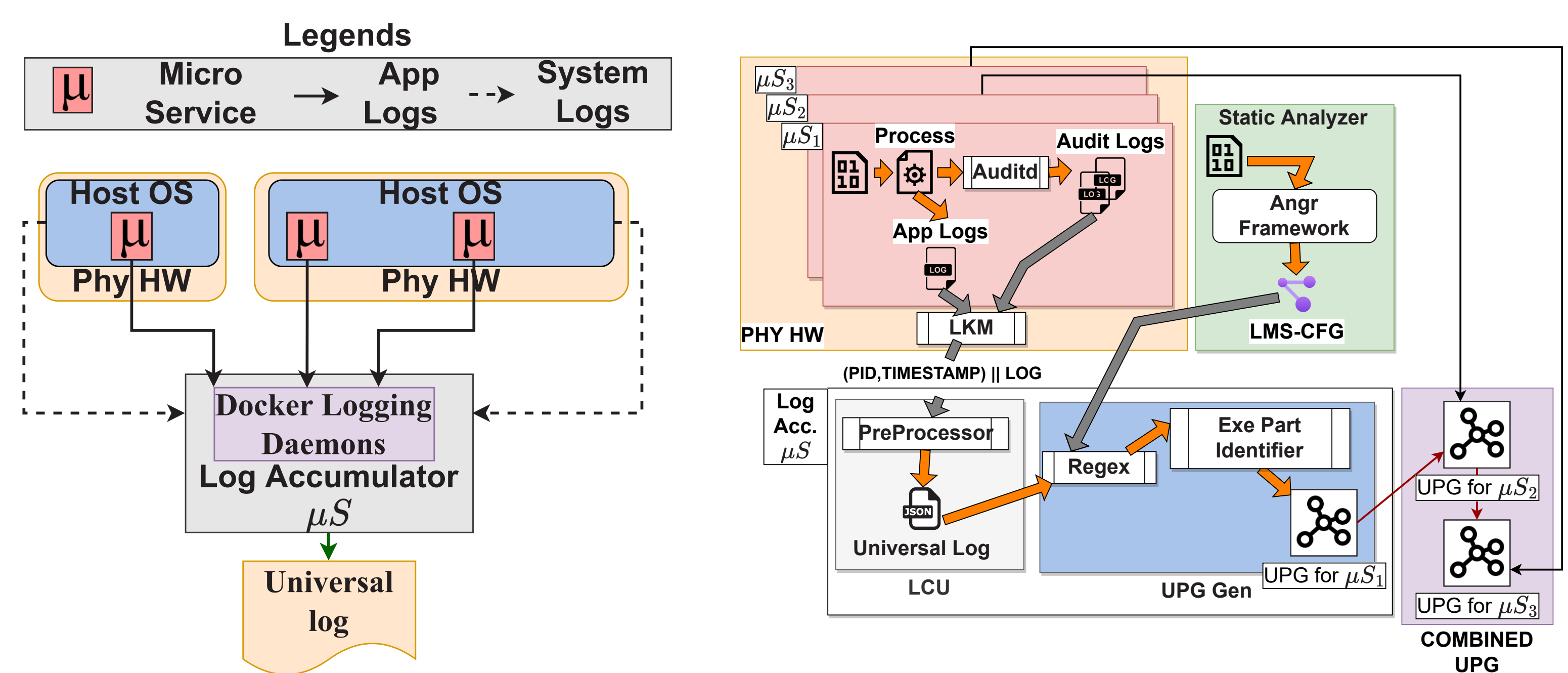


Figure 3. Left side, Environment Setup. Right side, DisProTrack Execution Flow.

## Results and Analysis

1) Since DisProTrack is targeted for serverless applications; resource overhead is a major concern. Therefore, experimentally we want to understand the resource overhead of the framework.

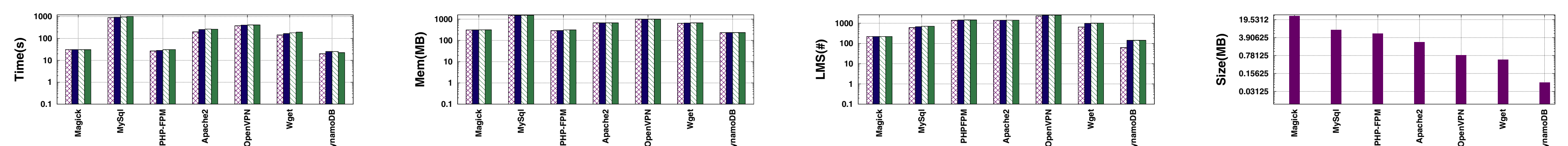


Figure 4. Performance Evaluation - Static Analysis - The Y-axes are in logarithmic scale

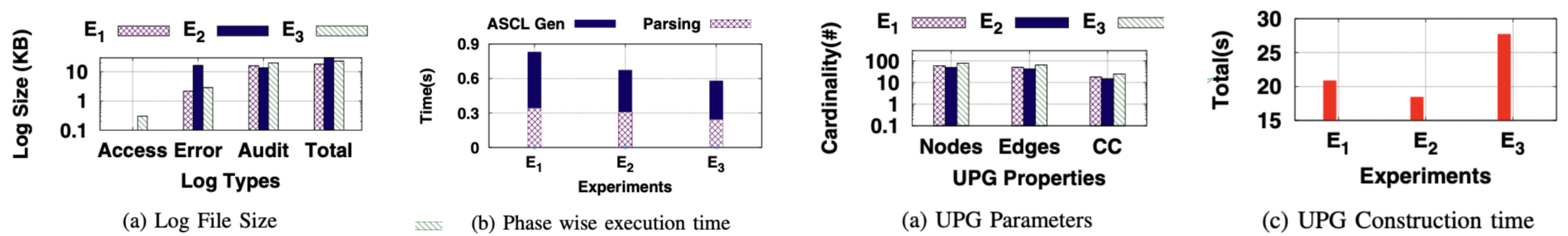


Figure 5. Performance Evaluation - Runtime Analysis

2) We also want to understand how effective DisProTrack is for identifying malicious activities

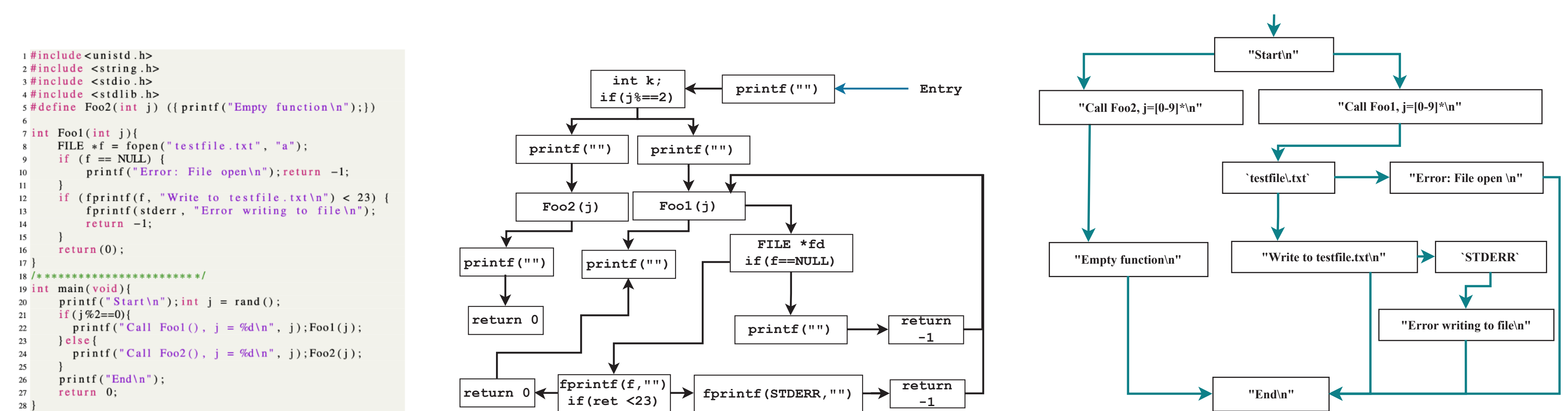


Figure 6. Left Side: PoC Program, Middle and Right Side: A PoC Case Study to Analyze the Accuracy of DisProTrack

## Conclusion

1. DisProTrack can be deployed as a microservice on top of the SLC without instrumenting the source code of the applications
2. Implementation is open-sourced
3. DisProTrack has a minimal memory footprint (~KB) & responds within 20s-30s.

### Acknowledgement

I would like to thank the "Infocom 2023" organizers for awarding me with the student travel grant