

# Engineering Report for the Scalable Compositional Static Taint Analysis for Sensitive Data Tracing on Industrial Micro-Services Research

Zexin Zhong

*University of Technology Sydney*  
Sydney, Australia  
*Ant Group*  
Hangzhou, China  
zhongzexin.zzx@antgroup.com

Jiangchao Liu

*Ant Group*  
Hangzhou, China  
jiangchao.ljc@antgroup.com

Diyu Wu

*Ant Group*  
Hangzhou, China  
wudiyu.wdy@antgroup.com

Peng Di

*Ant Group*  
Hangzhou, China  
dipeng.dp@antgroup.com

Yulei Sui

*University of Technology Sydney*  
Sydney, Australia  
yulei.sui@uts.edu.au

Alex X. Liu

*Ant Group*  
Hangzhou, China  
alexliu@antgroup.com

John C.S. Lui

*Chinese University of Hong Kong*  
HongKong, China  
cslui@cse.cuhk.edu.hk

***Index Terms*—program analysis, taint analysis, micro-services**

## I. ADAPTING TO ENTERPRISE-SPECIFIC FRAMEWORKS: A PRACTICE IN SOFA-BASED FINTECH MICRO-SERVICES

Different enterprises have different frameworks based on their business scenario and technology stack. In this section, we introduce how we resolve the framework problems in SOFA architecture. We first performed the modeling task on the target application before the compositional taint analysis. Then, the modeling result of the target applications will be processed by the taint analysis module for the analysis. This automated process includes preprocessing of the SOFA framework (the financial micro-services variant framework of Spring), RPC and database operation.

**Preprocessing and Modeling.** According to our experience, the AOP, Spring Bean injections, and SOFA Bean injections features that provided by the SpringBoot [1] and SOFA [2] framework usually lead to missing knowledge about the instance type of some object or the particular code that need to be executed at some function. AOP, Spring Beans, and SOFA Beans are usually configured in configuration XML files and are used in the form of annotation in the code at the beginning or end of the function. Thus, in our case, we will first automatically scan and extract the above information from the configuration XML files based on the semantics of the SOFA-RPC and Spring XML domain language. Then we model this information in a specified format by using the code transformation and modeling supporter. These models will be loaded into the memory for the further analysis task. In our compositional taint analysis process, when the specified annotation is found, the analyzer will query the model from

memory based on the annotation details and use the loaded model for the taint analysis.

**Entry Point Determination.** Because our proposed approach, which used the compositional analysis approach, does not rely on the traditional top-down analysis mode, therefore, determining the specified entry points is not required for our analysis. Therefore, distinct from other approaches based on the traditional top-down analysis mode, we do not need to identify the specified entry point before the analysis. The analyzer will randomly select an entry analysis function and execute the analysis process on all functions in the program to cover the whole program.

**Sources and Sinks Determination.** According to the successful experience of ANTaint [3] in sources and sink determination, we also assume that all the incoming data in the program will be marked as sources. These incoming data include the input parameters of the entry points, the return value of the RPC function, the interested fields or parameters provided by the developers, and the fields or objects that query from the database by the 'Select' operation. Furthermore, we assume that all the out-coming (output) data or the output functions which will expose the data to other environments or namespace will be marked as sinks. These out-coming (output) data include all the return values from the entry points, parameters for the RPC functions, the log or system printing functions, and the object or fields that will be updated to the database by the 'update', 'insert', and 'delete' operation.

## II. INDUSTRY APPLICATION

The introduced field-based compositional static taint analysis approach is used in various industry scenarios of tracing sensitive data for change governance, data breach detection, and data consistency validation at Ant Group.

**Change governance.** Consider there is a database that contains a transaction history table and the database will be queried by multiple services more than 10 million times per day. The primary key '*transaction\_id*' in the transaction history table is an 'int' type field, which will be auto increased for each new transaction record. With the increasing transaction volume, an 'int' type '*transaction\_id*' is not enough for future transactions (maximum 4,294,967,295 records in total). To address this issue, the database team needs to migrate the data to a new database table and update the data type of the '*transaction\_id*' from 'int' to 'BigInt'. In this case, developers need to notify those micro-services that will query data from this table with the '*transaction\_id*', to make corresponding updates. Developers need to update the data type of parameter related to the '*transaction\_id*' in the micro-services from 'int' to 'long' or 'BigInteger'. Otherwise, the overflow error will be triggered when loading the 'BigInt' type of data of the '*transaction\_id*' field from the database. Manually identifying all these related parameters and all data flows is a complicated task, which is time-consuming due to the scale and complexity of the industry micro-services applications. Furthermore, although manually identified is possible in some cases, the accuracy needs to be verified. For the site reliability requirement, to guarantee the reliability of the system, a verification mechanism needs to be used to verify these changes. Therefore, developers need to verify that all parameters related to the change are manually identified correctly and updated in the micro-services. In the above scenes, a scalable analysis tool that can trace the sensitive data or specified field in the large-scale complex micro-services is necessary for any change governance. As the analysis can help to identify and verify all parameters that are related to the change in the micro-services automatically.

## REFERENCES

- [1] R. Johnson, *Spring*, 2021. [Online]. Available: <https://spring.io/>
- [2] SOFAS*Stack*, *SOFAS<sub>Stack</sub>*, 2021. [Online]. Available: <https://www.sofastack.tech/en/>
- [3] J. Wang, Y. Wu, G. Zhou, Y. Yu, Z. Guo, and Y. Xiong, “Scaling static taint analysis to industrial soa applications: A case study at alibaba,” in *ESEC/FSE’20*, 2020, p. 1477–1486.