

# Linhai Song - Research Statement

My research interests span the areas of systems, security, and software engineering. The goal of my research is to help developers build more efficient, reliable, and secure software systems.

My dissertation research centers around software performance. Naturally, everyone wants software to run fast. Slow and inefficient software can easily frustrate end users and cause economic loss. The software-inefficiency problem has already caused several highly publicized failures. My research philosophy is to view the software-inefficiency problem from the perspective of combating *performance bugs*. Performance bugs are software implementation mistakes that can cause inefficient execution. Performance bugs cannot be optimized away by state-of-practice compilers. Many of them escape the in-house testing and manifest in front of end users, causing severe performance degradation and a huge waste of energy in the field. Performance bugs are becoming more critical, with the increasing complexity of modern software and workload, the meager increases of single-core hardware performance, and pressing energy concerns. It is urgent to combat performance bugs.

To fight performance bugs, my methodology is to investigate existing approaches originally designed for functional bugs and try to apply, adapt, and extend them for performance bugs. My experience spans different stages of combating performance bugs: real-world bug understanding, bug detection, failure diagnosis, and automated bug fixing. In particular, I conducted the first characteristics study on performance bugs, based on 110 bugs randomly collected from five representative open-source software suites [10]. I built a static performance bug detection tool suite by using extracted efficiency rules from fixed performance bugs [10], and helped build a dynamic performance bug detection tool for inefficient nested loops [9]. I explored how to apply statistical debugging to performance failure diagnosis [8], and designed a series of static-dynamic hybrid analysis routines to provide more detailed diagnosis information for inefficient loops [4]. I designed three source-to-source code transformations to automatically fix performance bugs in parallel applications, which offload computation to Intel Xeon Phi manycore coprocessors [7].

Besides software efficiency, I also worked on two software reliability projects as part of my doctoral research. In the first, I helped build a concurrency bug fixing system [11], which can automatically eliminate atomicity-violation bugs. In the second, I helped study change histories of critical sections in open-source software [6] to provide a better understanding of synchronization challenges faced by real-world developers.

After graduation, I started to work on security, with a focus on applying big data techniques to analyze and predict security incidents. One leveraged data repository is VirusTotal, which contains billions of real-world malware labeled by state-of-the-art anti-virus engines. I studied characteristics of real-world malware, modeled how influence propagates across different anti-virus vendors, and explored the feasibility of building machine learning malware detectors based only on hash values of files [3, 5]. The study results can shed light on future research directions and assist both anti-virus vendors and normal users in their fight against malware.

My research work has already had industrial and academic impact. My two performance bug detection techniques [9, 10] found hundreds of previously unknown performance bugs in mature open-source software, many of which have already been confirmed and fixed by developers. My performance bug fixing project [7] won MICRO'14 best paper runner up, and has also led to an issued patent [1]. My concurrency bug fixing system [11] won ACM SIGPLAN Research Highlights Award [2]. As of now, my publications have 415 citations in total.

## 1 Dissertation Research

To address the software-inefficiency problem, my research efforts cover the whole stack of software systems, from hardware to software, and cover all aspects of combating software bugs, from understanding, to detecting, diagnosing, and fixing.

**Real-world Performance Bug Understanding.** Like functional bugs, research on performance bugs should also be guided by empirical studies. Poor understanding of performance bugs is part of the causes of today's performance bug problem. In order to improve the understanding of real-world performance bugs, I co-led the first, to the best of our knowledge, empirical study on real-world performance bugs, based on 110 bugs randomly sampled from five open-source software suites [10]. Following the lifetime of performance bugs, our study was mainly performed in four dimensions. We studied the root causes of performance bugs, how they are introduced, how to expose them, and how to fix them. The main findings of our study include: (1) performance bugs have dominating root causes and fix strategies, which are highly correlated with each other; (2) workload mismatch and misunderstanding of APIs' performance features are two major reasons why performance bugs are introduced;

and (3) around half of the studied performance bugs require inputs with both special features and large scales to manifest. Our empirical study can guide future research on performance bugs. According to Google Scholar, our paper has already been cited 125 times since 2012. Our empirical study has already motivated our own performance bug detection and performance failure diagnosis projects.

**Static/Dynamic Performance Bug Detection.** Our empirical study shows that both statically checkable efficiency rules and violations of these rules exist widely in software. Inspired by this finding, we manually inspected final patches of fixed performance bugs in our studied performance bug set, extracted efficiency rules from 25 bug patches, and implemented static checkers to detect rules’ violations [10]. In total, our static checkers found 332 previously unknown performance bugs from the latest versions of Apache, Mozilla, and MySQL. We reported some of identified bugs to developers. 77 reported bugs have already been confirmed by developers, including 15 reported bugs already fixed by developers. Our empirical study also finds that 90% of performance bugs involve loops, and 50% of performance bugs involve at least two levels of loops. Motivated by this finding, I helped a fellow graduate student build a novel automated test oracle named Toddler [9], which enables testing for performance bugs caused by inefficient nested loops to use the well-established and automated process of testing for functional bugs. Using Toddler, we found 42 new bugs in six Java projects. Based on our bug reports, developers so far have fixed 21 bugs and confirmed 6 more as real bugs.

**Performance Failure Diagnosis.** Due to the preliminary tool support, many performance bugs escape the in-house performance testing and manifest in front of end users. After users report performance bugs, developers need to diagnose them and fix them. Diagnosing user-reported performance failure is another key aspect of fighting performance bugs.

We first investigated the feasibility and design space to apply statistical debugging to performance failure diagnosis [8]. After studying 65 user-reported performance bugs in our bug set, we found that the majority of performance bugs are observed through comparison, and many user-filed performance bug reports contain not only bad inputs, but also similar and good inputs. Statistical debugging is a natural fix for user-reported performance bugs. We evaluated three types of widely used predicates and two representative statistical models. Our evaluation results show that branch predicate plus two statistical models can effectively diagnose user-reported performance failure. The basic model can help diagnose performance failure caused by wrong branch decision, and the  $\Delta$ LDA model can identify inefficient loops. We applied sampling to performance failure diagnosis. Our experimental results show that some unique nature of loop-related performance bugs allows sampling to lower runtime overhead without sacrificing diagnosis latency, which is very different from functional failure diagnosis.

We then built LDoctor [4] to provide more fine-grained diagnosis information for inefficient loops through a two-step process. We first figured out a root-cause taxonomy for common inefficient loops through a comprehensive study on 45 inefficient loops. Our taxonomy contains two major categories: resultless and redundancy, as well as several subcategories. Guided by our taxonomy, we then designed a series of analysis for inefficient loops. Our analysis focuses its checking on suspicious loops pointed out by statistical debugging, hybridizes static and dynamic analysis to balance accuracy and performance, and relies on sampling and other designed optimization to further lower runtime overhead. Evaluation results under real-world inefficient loops show that LDoctor can cover most root-cause subcategories, report few false positives, and bring a low runtime overhead.

**Performance Bug Fixing.** Intel Xeon Phi coprocessors have recently been introduced as new members of the manycore family. Compared with GPU, Xeon Phi coprocessors are easier to program, since they provide x86 compatibility and support many different programming models. To offload existing parallel loops, developers just need to add simple pragmas. However, our recent experience shows that simply adding pragmas does not result in better performance, and too many performance bugs are contained in offloaded parallel loops.

After careful investigation, we designed three source-to-source code transformations to automatically fix performance bugs contained in offloaded parallel loops [7]. The first transformation, data streaming, is designed to overlap data transfers between CPUs and coprocessors with computation on coprocessors to hide data transfer overhead and reuse memory on coprocessors. The second transformation, regularization, is designed to re-arrange the order of computation and regularize loops with irregular memory accesses. The last transformation is designed to support efficient transfer for large pointer-based data structures between CPUs and coprocessors. The designed transformations can benefit 9 out of 12 benchmarks in our experiments, and improve performance by 1.16x - 52.21x. This work won MICRO’14 best paper runner up.

## 2 Future Research

Going forward, I hope to leverage my areas of expertise to improve the performance, reliability, and security of various types of computer systems.

**Performance.** My intermediate research goal is to combat performance bugs in different aspects from the work I have already done. For example, I would like to explore how to monitor algorithmic complexity during production runs. On-line algorithmic monitoring can help developers understand real-world workloads, identify codes in superlinear complexity, and expose new performance optimization opportunities. Previous works on algorithmic profiling or monitoring will incur more than 10x runtime overhead, which cannot be tolerated in production runs. I plan to build a tool to collect runtime information with a reasonable overhead and infer approximate complexity from deployed software. As another example, I would like to provide tool support for performance testing. Our empirical study shows that almost half of studied performance bugs need inputs with both special features and large scales to manifest. Existing techniques are designed to generate inputs with good code coverage and focus only on special features. I plan to extend existing input generation techniques with emphasis on large scales. Another important problem during performance testing is to automatically judge whether a performance bug has occurred. I plan to leverage existing dynamic performance bug detection techniques to build performance testing oracles.

**Reliability.** As big data is changing every single business, more and more developers are working on using big data computing systems, like Hadoop and Spark, to process their massive amounts of data. My own experience in leveraging Spark to analyze VirusTotal’s data has given me insight into the challenges of debugging big data applications. First, it is time-consuming to repeat a failure, and it may take developers several hours to receive notice of a runtime failure or an incorrect output. Second, it is almost infeasible to identify failure-triggering inputs among millions of input records without tool support. Third, it is difficult to identify a failure’s root cause, since the failure may propagate across different stages and across different nodes.

My longer term research direction is to provide tool support for debugging big data applications. I plan to start by studying real-world bugs in big data applications to figure out whether there are common bug patterns. Then I would like to build static bug detectors to identify bugs following common patterns before executing big data applications. To quickly identify failures’ root causes, I plan to build slicing tools that can analyze execution across different nodes, interactive breakpoint tools that can stop the execution for single records specified by developers, and log mining tools that can analyze large quantities of logs in big data applications. In order to patch bugs in big data applications, I would like to build automatic fixing tools. It is time-consuming to run the patched version from the very beginning, and I would like to build tools which can help reuse computation already finished in failure runs.

**Security.** I plan to continue my research on security along two directions. One is to apply big data techniques to security, and the other is to detect vulnerabilities in devices from the Internet of Things (IoT) ecosystem.

Big security data provides opportunities to leverage data-centric methods to improve today’s security techniques. For example, we can study what attackers have done on a large scale to figure out their strategies and predict what they will do in the future. As another example, by utilizing malware that has already been studied and labeled by security experts, we can learn how malware evolves and quickly detect emerging malware. We can also apply deep learning techniques on these labeled malware samples and build new malware detectors. It is also important to think about whether existing big data computing systems are suitable for security workload, and look for opportunities to improve current systems.

Vulnerability detection in the IoT ecosystem is more critical than ever before. As of now, there are billions of connected devices used worldwide, and this number is growing rapidly. Hacked devices can leak sensitive information and power denial of service attacks. Due to limited computation resources on a single device, anti-virus software cannot be used to prevent attacks. To make things worse, security is usually an afterthought for developers who build firmware for these devices. The need to detect vulnerabilities in firmware earlier and patch them earlier is increasing dramatically. To detect vulnerabilities, I plan to study mechanisms for various firmware and devices, and try to extend existing techniques, like symbolic execution and whitebox fuzzing. After identifying a new vulnerability, it is also important to quickly search affected devices. I plan to build an efficient bug search system for IoT products.

## References

- [1] Min Feng, Srimat Chakradhar, and **Linhai Song**. “Compiler Optimization for Many Integrated Core Processors”. U.S. Patent No. 20150277877.
- [2] SIGPLAN. “SIGPLAN Research Highlights Papers”. <http://www.sigplan.org/Newsletters/CACM/Papers/>.
- [3] **Linhai Song**, Ce Zhang, Yiyang Zhang, Heqing Huang, and Wu Zhou. “Understanding Malwares and Anti-Virus Engines in the Real World”. [http://songlh.github.io/eurosys\\_submission.pdf](http://songlh.github.io/eurosys_submission.pdf) (**Under Submission**).
- [4] **Linhai Song** and Shan Lu. “Performance Diagnosis for Inefficient Loops”. In *Proceedings of the 2017 International Conference on Software Engineering, ICSE’17*, May 2017.
- [5] **Linhai Song**, Heqing Huang, Wu Zhou, Wenfei Wu, and Yiyang Zhang. “Learning from Big Malwares”. In *Proceedings of the 7th ACM SIGOPS Asia-Pacific Workshop on Systems, APSys’16*, 2016.
- [6] Rui Gu, Guoliang Jin, **Linhai Song**, Linjie Zhu, and Shan Lu. “What Change History Tells Us About Thread Synchronization”. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering, ESEC/FSE’15*, 2015.
- [7] **Linhai Song**, Min Feng, Nishkam Ravi, Yi Yang, and Srimat Chakradhar. “COMP: Compiler Optimizations for Manycore Processors”. In *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture, MICRO’14*, December 2014. (**Best Paper Runner Up**).
- [8] **Linhai Song** and Shan Lu. “Statistical Debugging for Real-world Performance Problems”. In *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages and Applications, OOPSLA’14*, October 2014.
- [9] Adrian Nistor, **Linhai Song**, Darko Marinov, and Shan Lu. “Toddler: Detecting Performance Problems via Similar Memory-access Patterns”. In *Proceedings of the 2013 International Conference on Software Engineering, ICSE’13*, May 2013.
- [10] Guoliang Jin\*, **Linhai Song\***, Xiaoming Shi, Joel Scherpelz, and Shan Lu. “Understanding and Detecting Real-World Performance Bugs”. In *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI’12*, June 2012. \*: equal contributions.
- [11] Guoliang Jin, **Linhai Song**, Wei Zhang, Shan Lu, and Ben Liblit. “Automated Atomicity-Violation Fixing”. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI’11*, June 2011. (**SIGPLAN Research Highlights Award**).