

# Supplemental Material for the Paper: “Insights from Bugs in FPGA High-Level Synthesis Tools: An Empirical Study of Bambu Bugs”

Zun Wang, He Jiang, *Member, IEEE*, Xiaochen Li, Shikai Guo, Xu Zhao, and Yi Zhang

**Abstract**—High-Level Synthesis (HLS) tools have been widely used in FPGA design to convert C/C++ code to Hardware Description Language (HDL) code. Unfortunately, HLS tools are susceptible to bugs which can introduce serious vulnerabilities in FPGA products, leading to substantial losses. However, the characteristics of these bugs (e.g., root causes and bug-prone stages) have never been systematically studied, which significantly hinders developers from effectively handling HLS tool bugs. To this end, we conduct the first empirical study to uncover HLS tool bug characteristics. We collect 349 bugs of a widely used HLS tool, namely Bambu. We study the root causes, buggy stages, and bug fixes of these bugs by applying a multi-person collaboration method. Finally, 13 valuable findings are summarized. We find 14 categories of root causes in Bambu bugs; most bugs (22.1%) are caused by incorrect implementation of IR processing; the front end of Bambu is more bug-prone; to fix these bugs, 2.27 files and 80.19 lines of code need to be modified on average. We also present the insights gained from 95 Vitis HLS bugs. From these findings, we suggest that developers could use an on-the-fly code generator configuration method to generate suitable testing programs for HLS tool bug detection and apply Large Language Models to assist in fixing HLS tool bugs.

**Index Terms**—FPGA, high-level synthesis tool, software bugs, bug characteristics, empirical study.

## I. COMPARISON WITH SIMILAR STUDIES DONE ON SOFTWARE COMPILERS

IN the supplemental material, we have made a comparison with similar studies done on software compilers [1], [2], [3], [4], [5]. More specifically, Sun et al. [1] conducted the first systematic study to analyze GCC and LLVM bugs. However, they did not analyze the root causes, bug-prone stages, and bug symptom of compiler bugs. Zhou et al. [2] conducted the first study to characterize the optimization bugs

of GCC and LLVM. Wang et al. [3] and Romano et al. [4] focused on analyzing the bug in the WebAssembly framework (i.e., WebAssembly compilers and runtime). Shen et al. [5] performed an exhaustive analysis of bugs in deep learning compilers. The details of the comparison are shown in the following subsections. Based on the comparison, we have shown that our study does not have a similar taxonomy and bug distribution compared to these studies.

### A. Root causes comparison

The root causes and their distribution of each study are shown in Table I. We have highlighted the top 3 root causes reported in each study. We exclude the studies from Zhou et al. [2], since they did not analyze the root causes in their studies. Wang et al. [3] concluded 16 root causes of WebAssembly runtime bugs (e.g., incorrect algorithm implementation, API misuse, and incorrect assignment). They found that incorrect algorithm implementation, memory, and incorrect exception handling are the top 3 root causes. Romano et al. [4] found 8 root causes of WebAssembly compiler bugs, such as asyncify synchronous code, incompatible data types, memory model difference, and other infrastructures. The top 3 root causes are incompatible data types, other infrastructures, and emulating native environment. Shen et al. [5] found 12 root causes of Deep Learning compiler bugs (e.g., type problem, API misuse). They stated that type problem is the most common root cause. Compared to these studies, we have 14 categories of root causes. There may be some overlap between our study and other studies, such as type problems and configuration problems. However, most root causes we found do not appear in previous studies (e.g., *Incorrect Implementation of Binding Processing*). For the distribution of root causes, we found that *Incorrect Implementation of IR Processing*, *Configuration File and Script Problem*, and *Incorrect Implementation of Code and File Generation* are the top 3 root causes of HLS tool bugs. Thus, the distribution of root causes in our study is also different from those found in other studies.

### B. Bug-prone stages comparison

The bug-prone stages and their distribution of each study are shown in Table II. We exclude the studies from Zhou et al. [2] and Wang et al. [3], since they did not analyze the bug-prone stages in their studies. For Romano et al. [4], they found build time, program compile time, and runtime bugs in WebAssembly compilers. The program compile time and

He Jiang is with the School of Software, Dalian University of Technology, Dalian 116024, China, also with the Key Laboratory for Ubiquitous Network and Service Software of Liaoning Province, Dalian 116024, China, and also with Dalian Key Laboratory of Artificial Intelligence, Dalian 116024, China, and also with DUT Artificial Intelligence Institute, Dalian 116024, China (corresponding e-mail: jianghe@dlut.edu.cn).

Zun Wang and Xiaochen Li are with the School of Software, Dalian University of Technology, Dalian 116024, China, and also with the Key Laboratory for Ubiquitous Network and Service Software of Liaoning Province, Dalian 116024, China (e-mail: wangzun\_ssdt2020@mail.dlut.edu.cn; xiaochen.li@dlut.edu.cn;).

Shikai Guo is with the School of Information Science and Technology, Dalian Maritime University, Dalian 116026, China (e-mail: shikai.guo@dlnu.edu.cn).

Xu Zhao and Yi Zhang are with the School of Computer Science and Technology, Beijing Institute of Technology, Beijing 100081, China (e-mail: zhaoxu@bit.edu.cn; 3220195123@bit.edu.cn).

TABLE I  
ROOT CAUSES COMPARISON

Compilers	Root causes	Notes
WebAssembly runtimes [3]	1. <b><i>Incorrect Algorithm Implementation (25.49%)</i></b> .	
	2. <b><i>Memory (17.07%)</i></b> .	
	3. <b><i>Incorrect Exception Handling (10.96%)</i></b> .	
	4. API Misuse.	
	5. Type Problem.	
	6. Incorrect Condition Logic.	
	7. Incorrect Assignment.	
	8. Incorrect Configuration.	
	9. Missing Condition Checks.	
	10. Concurrency.	
	11. Environment Incompatibility.	
	12. External API Incompatibility.	
	13. Dependent Module Issue.	
	14. Logical Order Error.	
	15. Incorrect Numerical Computation.	
	16. Others.	
WebAssembly compilers [4]	1. Asyncify Synchronous Code.	The author did not present the top 3 root causes in the paper. The number behind the root cause represents the number of bugs related to that specific root cause. According to the number of bugs, we have highlighted the top 3 root causes.
	2. <b><i>Incompatible Data Types (23)</i></b> .	
	3. Memory Model Difference.	
	4. <b><i>Other Infrastructures (25)</i></b> .	
	5. <b><i>Emulating Native Environment (23)</i></b> .	
	6. Supporting Web APIs.	
	7. Cross-Language Optimizations.	
	8. Runtime Implementation Discrepancy.	
Deep Learning compilers [5]	1. API Misuse.	
	2. Incompatibility.	
	3. <b><i>Type Problem (19.23%)</i></b> .	
	4. <b><i>Tensor Shape Problem (13.27%)</i></b> .	
	5. Incorrect Numerical Computation.	
	6. Incorrect Assignment.	
	7. Incorrect Exception Handling.	
	8. Misconfiguration.	
	9. Concurrency.	
	10. <b><i>Incorrect Code Logic (16.92%)</i></b> .	
	11. Typo.	
	12. Others.	
HLS tools (our study)	1. Incorrect Compiler Usage.	
	2. Incorrect Implementation of Compiler Plug-ins.	
	3. Incorrect Implementation of Tool Parameter Processing.	
	4. Type Problem.	
	5. <b><i>Configuration File and Script Problem (14.6%)</i></b> .	
	6. Incorrect Implementation of Specific Structure Data Processing.	
	7. <b><i>Incorrect Implementation of IR Processing (22.1%)</i></b> .	
	8. Incorrect Implementation of Allocation Processing.	
	9. Incorrect Implementation of Scheduling Processing.	
	10. Incorrect Implementation of Binding Processing.	
	11. Incorrect Implementation of Architecture and Interface Processing.	
	12. <b><i>Incorrect Implementation of Code and File Generation (13.2%)</i></b> .	
	13. Incorrect Implementation of Tool Integration.	
	14. Other Problems.	

runtime are relatively more buggy than the build time. For Shen et al. [5], they found three bug-prone stages in deep learning compilers, including high-level IR transformation stage, low-level IR transformation stage, and model loading stage. They presented that the high-level IR transformation stage (44.92%) and the low-level IR transformation stage (33.64%) are relatively more buggy than the model loading stage (21.44%) on average. In our study, we found 5 bug-prone stages, such as build, configuration, front end, HLS process, and back end. We found the front end, HLS process, and back end are more bug-intensive stages in HLS tools. We think our analysis results are different from those reported in other studies due to differences in architecture between HLS tools and other compilers. For example, there is no HLS process in

deep learning compilers.

### C. Bug symptoms comparison

The bug symptoms and their distribution of each study are shown in Table III. Zhou et al. [2] found 3 bug types of compiler optimization bugs, including crash, misoptimization, and performance. They stated that misoptimization is the most common bug symptom in both GCC and LLVM, accounting for 57.21% and 61.38%, respectively. Wang et al. [3] et al. summarized 6 symptoms (i.e., crash, incorrect functionality, build error, bad performance, hang, and unreported). They found that crash is the most prevalent bug symptom. Romano et al. [4] found 11 types of bug symptoms and compile error is the most common bug symptom. Shen et al. [5] found 6

TABLE II  
BUG-PRONE STAGES COMPARISON

Compilers	Bug-prone stages	Notes
WebAssembly compilers [4]	1. <b><i>Program runtime (299)</i></b> . 2. <b><i>Program compile time (301)</i></b> . 3. Build time (95).	The author did not present the top 3 bug-prone stages in the paper. The number behind the stage represents the number of bugs related to that specific stage. According to the number of bugs, we have highlighted the top 3 stages.
Deep Learning compilers [5]	1. <b><i>High-level IR transformation stage (44.92%)</i></b> . 2. <b><i>Low-level IR transformation stage (33.64%)</i></b> . 3. Model loading stage (21.44%).	The author only presented these 3 stages in the paper.
HLS tools (our study)	1. Build. 2. Configuration. 3. <b><i>Front end (39.3%)</i></b> . 4. <b><i>HLS process (19.8%)</i></b> . 5. <b><i>Back end (20.3%)</i></b> .	

TABLE III  
BUG SYMPTOMS COMPARISON

Compilers	Bug symptoms	Notes
Optimization bugs in GCC and LLVM [2]	1. <b><i>Crash (39.33% and 34.08%)</i></b> . 2. <b><i>Misoptimization (57.21% and 61.38%)</i></b> . 3. Performance.	The author only presented these 3 bug symptoms (types) in the paper.
WebAssembly runtimes [3]	1. <b><i>Crash (56.86%)</i></b> . 2. <b><i>Incorrect Functionality (22.15%)</i></b> . 3. <b><i>Build Error (12.92%)</i></b> . 4. Bad Performance. 5. Hang. 6. Unreported.	
WebAssembly compilers [4]	1. <b><i>Build Error (95)</i></b> . 2. <b><i>Compile Error (264)</i></b> . 3. Linker Error. 4. Code Bloating. 5. Crash. 6. Data Corruption. 7. Fail to Instantiate. 8. Performance Drop. 9. Hang. 10. <b><i>Other Runtime Error (117)</i></b> . 11. Incorrect Functionality.	The author did not present the top 3 bug symptoms (types) in the paper. The number behind the bug symptom (type) represents the number of bugs related to that specific bug symptoms (types). According to the number of bugs, we have highlighted the top 3 bug symptoms (types).
Deep Learning compilers [5]	1. <b><i>Crash (59.37%)</i></b> . 2. <b><i>Wrong Code (25.04%)</i></b> . 3. Bad Performance. 4. Hang. 5. <b><i>Build Failure (8.29%)</i></b> . 6. Unreported.	
HLS tools (our study)	1. <b><i>Incorrect Results (38.9%)</i></b> . 2. <b><i>Tool Failure (33.7%)</i></b> . 3. <b><i>Incorrect GUI (9.5%)</i></b> . 4. Crash. 5. Hang. 6. Unreported.	

symptoms (i.e., crash, wrong code, bad performance, hang, build failure, and unreported) and crash is the most common bug symptom. In our study, we also found 6 types of bug symptoms, including incorrect results, tool failure, incorrect GUI, crash, hang, and unreported. Crash, hang, and unreported are common bug symptoms reported in many studies, as well in ours, but we also found tool failure, that is rarely mentioned in other studies. Meanwhile, we found that the symptom ‘incorrect results’ is the most common bug symptom in HLS tools, which is different from other studies where crash is the most common bug symptom.

## REFERENCES

- [1] C. Sun, V. Le, Q. Zhang, and Z. Su, “Toward understanding compiler bugs in gcc and llvm,” in *Proc. 25th Int. Symp. Softw. Test. Anal.*, 2016, pp. 294–305.
- [2] Z. Zhou, Z. Ren, G. Gao, and H. Jiang, “An empirical study of optimization bugs in gcc and llvm,” *J. Syst. Softw.*, vol. 174, pp. 110 884–110 884, 2021.
- [3] Y. Wang, Z. Zhou, Z. Ren, D. Liu, and H. Jiang, “A comprehensive study of webassembly runtime bugs,” in *Proc. IEEE Int. Conf. Softw. Anal., Evol. Reengineering*. IEEE, 2023, pp. 355–366.
- [4] A. Romano, X. Liu, Y. Kwon, and W. Wang, “An empirical study of bugs in webassembly compilers,” in *Proc. 36th IEEE/ACM Int. Conf. Autom. Softw. Eng.*. IEEE, 2021, pp. 42–54.
- [5] Q. Shen, H. Ma, J. Chen, Y. Tian, S.-C. Cheung, and X. Chen, “A comprehensive study of deep learning compiler bugs,” in *Proc. 29th ACM Jt. Meet. Eur. Softw. Eng. Conf. Symp. Found. Softw. Eng.*, 2021, pp. 968–980.