

# Native Code Level Test of Optimizing Performance of Android Compilers

Naoki YOSHIDA <sup>†,\*</sup> Toya HAMADA <sup>†</sup> Nagisa ISHIURA <sup>††</sup>

<sup>†</sup> Graduate School of Science and Technology <sup>††</sup> School of Engineering  
Kwansei Gakuin University

1 Gakuen Uegahara, Sanda, Hyogo, 669-1330, JAPAN

**Abstract—** In this paper, we introduce a technique for assessing the optimization performance of the Android DEX compiler at the level of native code. This method is designed to detect missed optimization in native codes generated by the Android runtime environment through random generation of Java programs. The detection of optimization deficiencies is performed using both differential method and equivalent program methods. In the differential method tests, we attempt to identify missed optimization by comparing the native code produced by newer and older versions of DEX compilers. In the equivalent program method, we aim to identify missed optimizations by comparing the native codes generated by a DEX compiler from both optimized and unoptimized source programs. The Random Java programs are generated from a modified version of Orange4, which were originally developed for generating C programs. The test systems, employing the proposed methods, effectively identified insufficient optimization in x86\_64 native code generated by the d8 DEX compiler.

## I. INTRODUCTION

Android<sup>1</sup> is an open-source operating system designed for mobile devices, and its utilization has been steadily on the rise in recent years. Android finds extensive application across various portable devices, including smartphones, which underscores the pressing need for improving its performance.

To tackle this challenge, continuous enhancements are being made to Android's compilers and runtime systems. As an alternative to the long-standing Dalvik virtual machine, Android adopted ART (Android Runtime) [1]. ART is based on Ahead-Of-Time (AOT) compilation, which precompiles parts of the code into native code. Consequently, it delivers superior runtime performance compared to Dalvik.

Alongside these developments, ensuring the reliability of these components through comprehensive testing is also a crucial task. In the context of testing Android's runtime, [2, 3, 4] have conducted testing of ART by generating random bytecodes. However, these tests primarily focus on error detection within the runtime systems and do not specifically address performance testing.

In contrast, this paper introduces a novel approach for evaluating the performance of Android's compilers and runtime. Our method aims to identify optimization shortcomings in both the DEX compiler and the AOT compiler in Android Runtime

(ART). Identification of optimization deficiencies is carried out based on both differential testing [5, 6] and equivalence testing [7]. Class files compiled from randomly generated Java programs are used as test inputs. We detect the missed optimization by analyzing the native code generated by the AOT compiler.

The experimental results with two test systems based on the proposed method have identified optimization deficiencies in the x86\_64 native code produced by the DEX compiler d8 and ART.

## II. RANDOM TESTING OF ANDROID PROCESSORS AND COMPILER OPTIMIZATIONS

### A. DEX compiler and Android virtual machine ART

The Android processing system is depicted in Fig. 1. First, a program undergoes compilation using either the Java or Kotlin compiler, resulting in the creation of a class file. Subsequently, the DEX compiler (dx or d8) processes the class file, producing a DEX file that encapsulates bytecode for Android's runtime environment, referred to as ART [1].

The virtual machine (VM) in ART executes bytecode through interpretation, but for performance improvement, a portion of the code is converted into native code by a backend compiler and executed directly by the processor. This compilation process is performed ahead of time (AOT) during app installation, not Just-In-Time (JIT).

Both DEX compilers, dx and d8, perform bytecode-level optimizations. Among them, d8 is the latest compiler, and it is known to exhibit equivalent or better runtime performance compared to dx. Additionally, d8 has shorter compilation times and generates smaller DEX files compared to dx.

ART takes the received DEX files and converts them into executable ELF format files called OAT files using dex2oat. Before compilation, dex2oat verifies the validity of the DEX files through verification and then extracts portions for compilation from bytecode through filtering.

In other words, dex2oat compiles only those methods that meet specific conditions, such as a certain method size (e.g., the number of instructions or maximum register numbers).

The backend compiler transforms the bytecode's selected portions into intermediate representation (IR) for optimization, and from IR, it generates native code. Finally, dex2oat adds the native code to the original DEX file's contents to create the OAT file.

\*Currently with NSW Inc., Japan.

<sup>1</sup><https://source.android.com> (accessed 2023-9-20)

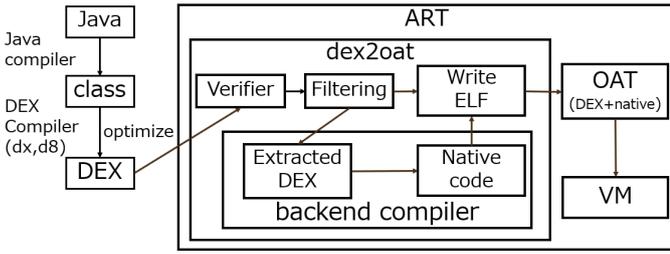


Fig. 1. DEX compiler and Android run time environment [4]

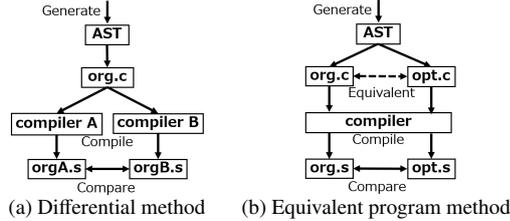


Fig. 2. Methods for testing compiler optimization [5]

```

01: #include <stdio.h>
02: #include <stdarg.h>
03: #define OK() printf("OK@\n")
04: #define NG(test,fmt,val) printf("@NG (test = " fmt ")@n",val)
05:
06: static volatile unsigned short x86 = 3U;
07: const signed short x88 = 386;
08: static volatile unsigned short t0 = 2428U;
09: volatile unsigned long long t1 = 2LLU;
10: static volatile signed char t3 = 9;
11: volatile signed char x92 = 115;
12:
13: int main (void)
14: {
15:     static const unsigned char x33 = 1U;
16:     static const volatile signed long long x41 = -1LL;
17:     volatile unsigned char x52 = 5U;
18:     const volatile unsigned char x65 = 38U;
19:     static signed long long x77 = 0LL;
20:     volatile signed int x87 = 511;
21:     const volatile signed long x89 = 2029215198009L;
22:     const volatile unsigned int x90 = 140U;
23:     volatile unsigned short t2 = 6570U;
24:     const unsigned long long x91 = 1255744988028LLU;
25:     static volatile unsigned long t4 = 1404136965112782LU;
26:     static volatile signed char t5 = 17;
27:
28:     t0 = ((unsigned short)(((signed int)x87)&((signed short)x88)));
29:     t1 = ((unsigned long)(((signed long)x89)*((signed long long)x33)));
30:     t2 = ((unsigned short)(((signed int)x90)<<(((signed char)x77)))));
31:     if(((signed char)((unsigned short)x86)*((signed int)x41))) {
32:         t3 = ((signed char)(((signed long)x77)<<(((signed long)x52)))));
33:         t4 = ((unsigned long)(((signed char)t3)|((unsigned long)x91)));
34:     }
35:     t5 = ((signed char)((signed char)x92)*((signed char)x65));
36:
37:     if (t0 == 386U) { OK(); } else { NG("t0", "%hu", t0); }
38:     if (t1 == 2029215198009LLU) { OK(); } else { NG("t1", "%llu", t1); }
39:     if (t2 == 140U) { OK(); } else { NG("t2", "%hu", t2); }
40:     if (t3 == 0) { OK(); } else { NG("t3", "%hhd", t3); }
41:     if (t4 == 1255744988028LU) { OK(); } else { NG("t4", "%lu", t4); }
42:     if (t5 == 1) { OK(); } else { NG("t5", "%hhd", t5); }
43:
44:     return 0;
45: }

```

Fig. 3. Test program generated by Orange4 [9]

### B. Random testing of compiler optimization

There are two distinct methods to assess the optimization performance of compilers: the differential method [5, 6] and the equivalent program method [7].

The differential method attempts to find optimization deficiencies in compilers by contrasting the outputs of two different compilers. The test procedure is illustrated in Fig. 2 (a). A test program `org.c` is compiled by distinct compilers (or different versions of the same compiler), `compilerA` and `compilerB`. The generated assembly codes denoted as `orgA.s` and `orgB.s` are then compared to identify optimization shortcomings, if any, in either of the compilers.

On the other hand, the equivalent program method tests a compiler using a pair of test programs which are functionally equivalent. In Fig. 2 (b), `org.c` and `opt.c` are the test programs, where `opt.c` is derived from `org.c` by applying anticipated optimization at the source code level. By comparing the assembly codes `org.s` and `opt.s` generated from them, optimization shortcomings are identified.

In terms of test programs for performance assessment, it is common to use benchmark programs or manually crafted test programs. However, for comprehensive testing, it is essential to utilize a large volume of test programs obtained through automatic generation.

In the context of automatic random generation of programs for compiler testing, a crucial challenge is how to avoid undefined behavior, such as division by zero, overflow, and infinite loops, all of which are tied to the dynamic behavior of these programs.

In the case of C program generation, tools like Csmith [8] achieve this by inserting guards to avoid undefined behavior. This approach facilitates the creation of programs that span a broad spectrum of C language syntax, but the forms of expressions or statements are limited.

In contrast, Orange4, as described in [9], guarantees the

elimination of undefined behavior through the utilization of data structures that consistently track values for all variables and expressions within the generated program. While this approach may impose certain restrictions on program syntax, it empowers Orange4 to generate a wide array of expressions, characterized by both diversity and length.

Fig. 3 provides an example of a test program generated by Orange4. Lines 6–26 declare variables, lines 28–35 contain statements with arithmetic expressions and control flow, and lines 37–42 involve statements for checking computed results.

## III. RANDOM TESTING OF ANDROID COMPILER OPTIMIZATION BASED ON NATIVE CODE COMPARISON

### A. Outline

This paper proposes novel random testing techniques designed to assess the optimization performance of the Android DEX compiler. The evaluation of the DEX compiler’s optimization performance is based on the quality of the native code generated within `dex2oat`. Please note that this evaluation takes place at the native code level, as the ultimate performance outcome is dictated by the execution of native code within the runtime environment, rather than the bytecode generated by the DEX compiler.

We present two distinct methods for performance testing, which are based on the differential method and the equivalent program method. The differential method involves comparing the native code outputs derived from the same Java program using two different DEX compilers. On the other hand, the equivalent program method revolves around the comparison of native code generated by a single compiler from two Java programs that are functionally equivalent.

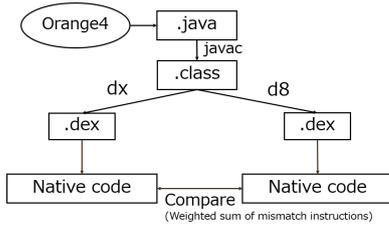


Fig. 4. Proposed flow of differential testing

The Java programs for the tests are generated by a modified version of Orange4 tailored for Java. Care is taken to generate only those those Java programs that meet the conditions for native code generation by dex2oat.

### B. Differential method

In the differential testing approach, two versions of DEX compilers are used. Typically, performance of a new version is tested against an old version, where the goal is to check whether the code generated by the newer version is not inferior to that generated by the older one.

The flow of the differential testing is illustrated in Fig. 4. Randomly generated Java programs are compiled, resulting in class files, which are then compiled using the two versions of DEX compilers, denoted as “dx” and “d8” in the figure, to generate two DEX files. Native codes generated by ART from each of these DEX files are extracted and compared to detect performance differences.

### C. Equivalent program method

In the equivalent program method, it is tested if a combination of a DEX compiler and AOT compilers perform intended optimization to generate native codes. The workflow of this method is illustrated in Fig. 5. In this approach, equivalent programs are generated by applying the expected optimizations at the source code level to randomly generated test programs. By compiling these programs using the d8 compiler and comparing the resulting two native codes, the goal is to detect optimization deficiencies.

In this paper, our specific focus lies in arithmetic optimization, which simplifies expressions during compile time. While this optimization is considered fundamental, it frequently exposes compiler bugs. Fig. 6 provides an illustrative example of arithmetic optimization in action. Initial constant values are assigned to the variables  $x_3$ ,  $x_4$ ,  $x_5$ , and  $x_6$ , and subexpressions are evaluated during the compilation process. Note that variables that variables  $x_1$  and  $x_2$  cannot be substituted with constants due to their declaration as `volatile`, indicating the potential for external updates to their values during program execution.

### D. Native code comparison

The native code segments to be compared are derived from the OAT files generated by ART from DEX files. A textual

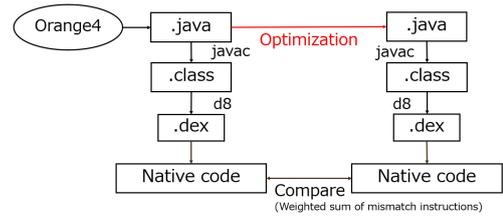


Fig. 5. Proposed flow of equivalence testing

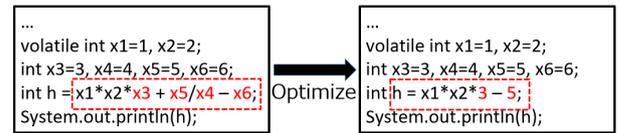


Fig. 6. Arithmetic optimization

dump of an OAT file encompasses various pieces of information, from which only the lines containing instructions and operands relevant to the target processor are isolated. An example of the extracted native code is illustrated in Fig. 7.

The comparison of native code is performed following the method described in [5], as shown in Fig. 8.

Initially, code pairs are examined from the start to the end, and contiguous sections where two codes exhibit matching instructions and operands for a minimum of  $k$  instructions (with  $k$  typically around 7) are designated as matching sections, while the remaining segments are classified as non-matching sections. Within a non-matching section, pairs featuring identical instructions are filtered out. The sum of the weights assigned to the remaining instructions (with higher weights assigned to instructions presumed to consume more execution cycles, such as multiplication/division and branch instructions) is assessed, and if the ratio surpasses a predefined threshold, it signifies a significant performance disparity in the code.

### E. Generation of random Java programs

We adapted the random C program generator of Orange4 for generating Java programs. Major modifications include the following three points.

#### 1. Variable modifiers

The generator has been adjusted to generate variable modifiers in accordance with the Java language specification. The modifiers are `volatile`, `final`, `static`, `private`, `protected`, `public`, and `transient`. Specifically, the `final` modifier is applied to variables that undergo a single assignment, while the `static` modifier is assigned to all class variables but not to local variables.

#### 2. Statements

The statements generated in the test programs are `if`, `for`, and `switch`. `while` statements are omitted due to a minor technical issue with Orange4.

#### 3. Explicit type conversion between integer and boolean types



TABLE II  
RESULT BY DIFFERENTIAL TESTING

(a) x86_64			(b) ARM		
target→ reference ↓	dx	d8	target→ reference ↓	dx	d8
dx	—	7	dx	—	0
d8	62	—	d8	13	—
3,000 tests			3,000 tests		

test1.java	
<pre> 1: class test 2: { 3:   static int [][] t0 = {{1}}; 4:   static int x1 = 1; 5: 6:   public static void main (String args[]) 7:   { 8:     int t1 = 1; 9:     t1 = 1/t0[0][0&amp;(1/x1)]; 10:  } 11: }</pre>	
test1.s (dx)	test1.s (d8)
<pre> ... 1: mov eax, [edx + 176] 2: test eax, eax 3: jz/eq +73 (0x000010ca) 4: 5: 6: 7: 8: 9: cmp [ebx + 8], 0 10: jbe/na +70 (0x000010d1) 11: mov eax, [ebx + 12] 12: test eax, eax 13: jz/eq +71 (0x000010dd) 14: add esp, 28 15: ret ...</pre>	<pre> ... 1: mov ecx, [ecx + 176] 2: test ecx, ecx 3: jz/eq +90 (0x000010db) 4: mov eax, 1 5: cmp ecx, -1 6: jz/eq +83 (0x000010e2) 7: cdq 8: idiv edx:eax, edx:eax / ecx 9: cmp [ebx + 8], 0 10: jbe/na +74 (0x000010e6) 11: mov eax, [ebx + 12] 12: test eax, eax 13: jz/eq +75 (0x000010f2) 14: add esp, 12 15: ret ...</pre>

Fig. 10. Test program detecting under-optimization of d8

test2.java	
<pre> 1: class test 2: { 3:   static int x0 = 1; 4: 5:   public static void main (String args[]) 6:   { 7:     int t0 = 1; 8:     int []x1 = {1} ; 9:     long [][]x2 = {{1,1}} ; 10:    t0 = (int)(x2[0&amp;(int)(x0*x2[0][0])][(1/x1[0])]); 11:  } 12: }</pre>	
test2.s (dx)	test2.s (d8)
<pre> ... 1: mov ecx, [ecx + 172] 2: xor eax, ecx 3: xor ecx, eax 4: xor eax, ecx 5: cdq 6: mov eax, [ecx + 8] 7: mov edx, [ecx + 16] 8: mov ebx, [ecx + 20] 9: cmp [ebx + 8], 0 10: 11: 12: 13: 14: 15: 16: 17: 18: 19: 20: 21: 22: 23: mov eax, 1 24: mov eax, eax 25: jnb/ae/nc +73 (0x00001114) 26: add esp, 20 27: pop ebp 28: pop esi 29: 30: ret ...</pre>	<pre> ... 1: mov ecx, [ecx + 172] 2: xor eax, ecx 3: xor ecx, eax 4: xor eax, ecx 5: cdq 6: mov ebx, [ecx + 8] 7: mov esi, [ecx + 16] 8: mov edi, [ecx + 20] 9: cmp [ebx + 8], 0 10: mov [esp + 4], esi 11: mov [esp + 8], edi 12: mov esi, eax 13: mov edi, edx 14: mov eax, [esp + 8] 15: imul eax, esi 16: imul edi, [esp + 4] 17: add edi, eax 18: mov eax, esi 19: mul edx:eax, eax * [esp + 4] 20: add edi, edx 21: mov esi, eax 22: mov eax, esi 23: mov eax, 1 24: mov eax, eax 25: jnb/ae/nc +77 (0x0000113d) 26: add esp, 32 27: pop ebp 28: pop esi 29: pop edi 30: ret ...</pre>

Fig. 11. Test program detecting under-optimization of d8

version of the ART subjected to testing was 2.1.0.

### A. Differential method

The results of the differential testing are presented in Table II. The target ISAs for testing were x86\_64 and ARM. In total, we executed 3,000 test programs, with the testing duration being approximately 3 hours for x86\_64 and around 6 hours for ARM.

In the table, “target” represents the compiler under evaluation, while “reference” refers to the compiler serving as a benchmark. The numerical values represent the quantity of programs for which the target compiler exhibited optimization shortcomings in comparison to the reference compiler. For example, in Table II (a), 62 programs identified optimization shortcomings in dx when contrasted with d8, and conversely, 7 programs identified optimization deficiencies in d8 when compared to dx.

Assuming that d8 is newer than dx, it is reasonable that dx exhibits higher count of optimization shortcomings than d8 in both (a) and (b). No optimization deficiencies were detected in d8 for ARM. On the other hand, 7 programs in x86\_64 detected optimization deficiencies in d8, indicating a regression in performance during the version upgrade. The 7 programs could be classified into two types.

A program of the the first type has been minimized, which

TABLE III  
RESULT BY EQUIVALENCE TESTING

Test programs	Under-optimization	Time [s]
3,000	558	13,229

is shown in Fig. 10. “test1.java” is a test program in Java and “test1.s (dx)” and “test1.s (d8)” are native codes generated by dx and d8, respectively. The native codes are identical except for lines 4–8, where d8 introduced generated an instruction sequence including division operation (idiv). The results of the idiv instruction are stored in the registers edx and eax, but they remain unused subsequently: at line 11 eax is reassigned another value at line 11 (it is also not used in the jump destination at line 10). This observation implies that d8 has failed to eliminate unnecessary codes in lines 4–8.

A minimized test program of the second type is shown in Fig. 11. The sequence of instructions d8 generated in lines 10–22 including multiplication operations (imul and mul) has no effect to the subsequent instructions and should have been removed through optimization.

org.java	opt.java
1: class test {	1: class test {
2: static volatile int x1 = 1;	2: static volatile int x1 = 1;
3: public static void	3: public static void
4: main(String args[]){	4: main(String args[]){
5: int x2 = 1;	5: int x2 = 1;
6: int x3 = 0;	6: int x3 = 0;
7: int h = (x3<<x2)*(1/x1);	7: int h = 0*(1/x1);
8: System.out.println(h);	8: System.out.println(h);
9: }	9: }
10: }	10: }
org.s	opt.s
1: jz/eq +92 (0x000010bc)	1: jz/eq +75 (0x000010ab)
2: mov eax, 1	2: mov eax, [RIP + 0xf9a]
3: cmp ecx, -1	3: cmpb [rax + 115], -16
4: jz/eq +86 (0x000010c4)	4: 4:
5: cdg	5: 5:
6: idiv edx:eax, edx:eax / ecx	6: 6:
7: mov eax, [RIP + 0xf89]	7: 7:
8: cmpb [rax + 115], -16	8: 8:
9: jb/nae/c +71 (0x000010c8)	9: jb/nae/c + 67 (0x000010b3)
10: test [rax + 7], 16	10: test [rax + 7], 16
11: mov esi, [rax + 244]	11: mov esi, [rax + 244]
...	...

Fig. 12. Test programs detecting under-optimization of d8

org.dex	opt.dex
...	...
0x0000: nop	0x0000: nop
0x0001: nop	0x0001: nop
0x0002: sget v1, I test.x1 // field@1	0x0002: sget v1, I test.x1 // field@1
0x0004: const/4 v0, #+1	0x0004: const/4 v0, #+1
0x0005: div-int/2addr v0, v1	0x0005: div-int/2addr v0, v1
0x0006: const/4 v1, #+0	0x0006: mul-int/lit8 v0, v0, #+0
0x0007: mul-int v1, v1, v0	0x0008: sget-object v1, Ljava/io/PrintStream;
0x0009: sget-object v0, Ljava/io/PrintStream;	java.lang.System.out // field@0
java.lang.System.out // field@0	0x000a: invoke-virtual [v0, v1], void
0x000b: invoke-virtual [v0, v1], void	java.io.PrintStream.println(int) // method@0
java.io.PrintStream.println(int) // method@0	0x000d: return-void
0x000c: return-void	...
...	...

Fig. 13. DEX codes for org.java and opt.java

## B. Equivalent program method

The d8 compiler was tested by 3,000 pair of programs, targeting the x86.64 ISA. The result is summarized in Table III. By this test, 558 programs identified optimization shortcomings in d8.

Fig. 12 shows one of the (minimized) test programs. “org.java” and “opt.java” are unoptimized and optimized Java programs, respectively, and “org.s” and “opt.s” are resulting native codes. Lines 4–8 of “org.s” which includes a division (idiv) instruction should be optimized away.

For reference, Fig. 13 is a part of the bytecode (DEX code), which is generated before being converted into native code by the AOT compiler in ART. At this stage, both programs generate the div instruction, indicating that the deficiency in optimization is not in the DEX compiler but rather in the AOT compiler of ART.

## V. CONCLUSION

This paper has proposed performance optimization test methods for Android compilers based on native code comparison. Through experiments, the both differential testing and equivalence testing successfully detect optimization deficiencies in the d8 compiler and the AOT compiler in ART.

Our current testing system is limited in its scope, primarily addressing a small portion of compiler optimization, specifically focusing on arithmetic optimization. Extending our test-

ing methodology to encompass other aspects of compiler optimization is a direction for future research.

## Acknowledgments

Authors would like to express their appreciation to the members of Ishiura Lab. of Kwansai Gakuin Univ. for their cooperation.

## REFERENCES

- [1] Android Open Source Project (Android core technologies), <https://source.android.com/devices/tech/dalvik/> (accessed 2021-12-17).
- [2] S. Kyle, H. Leather, B. Franke, D. Butcher, S. Monteith: “Application of Domain-aware Binary Fuzzing to Aid Android Virtual Machine Testing,” *ACM SIGPLAN Notices*, vol. 50, Issue 7, pp. 121132 (July 2015).
- [3] H. Ikeo, R. Shimizu, and N. Ishiura: “Random Testing of Android Virtual Machine by Valid DEX File Generation” (in Japanese), in *Technical Report of IEICE*, VLD2017-88, (Feb2018).
- [4] R. Shimizu and N. Ishiura: “Reinforcing Generation of Instruction Sequences in Random Testing of Android Virtual Machine” (in Japanese), in *Technical Report of IEICE*, VLD2018-124, (Mar. 2019).
- [5] K. Kitaura and N. Ishiura: “Random Testing of Compilers’ Performance Based on Mixed Static and Dynamic Code Comparison,” in *Proc. ACM International Workshop on Automating TEST Case Design, Selection, and Evaluation (A-TEST 2018)*, pp. 38–44 (Nov. 2018).
- [6] G. Barany: “Finding Missed Compiler Optimizations by Differential Testing,” in *Proc. International Conference on Compiler Construction (CC 2018)*, pp. 92–91 (Feb. 2018).
- [7] A. Hashimoto and N. Ishiura: “Detecting Arithmetic Optimization Opportunities for C Compilers by Randomly Generated Equivalent Programs,” *IPSJ Trans. System LSI Design Methodology*, vol. 9, pp. 21–29 (Feb. 2016).
- [8] X. Yang, Y. Chen, E. Eide and J. Regehr: “Finding and Understanding Bugs in C Compilers,” in *Proc. ACM Conference on Programming Language Design and Implementation (PLDI ’11)*, pp. 283–294 (Oct. 2011).
- [9] K. Nakamura and N. Ishiura: “Random Testing of C Compilers Based on Test Program Generation by Equivalence Transformation” in *Proc. Asia and Pacific Conference on Circuits and Systems (APCCAS 2016)*, pp. 676679 (Oct.2016).