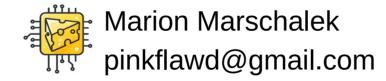
Compiler InternalsFor Security Engineers



Reflections on Trusting Trust

To what extent should one trust a statement that a program is free of Trojan horses? Perhaps it is more important to trust the people who wrote the software.

KEN THOMPSON

INTRODUCTION

What security relevant aspects do compilers provide?

- Mitigations
- Optimizations
- Obfuscation
- Analysis
- Instrumentation
- Intermediate representations to no end

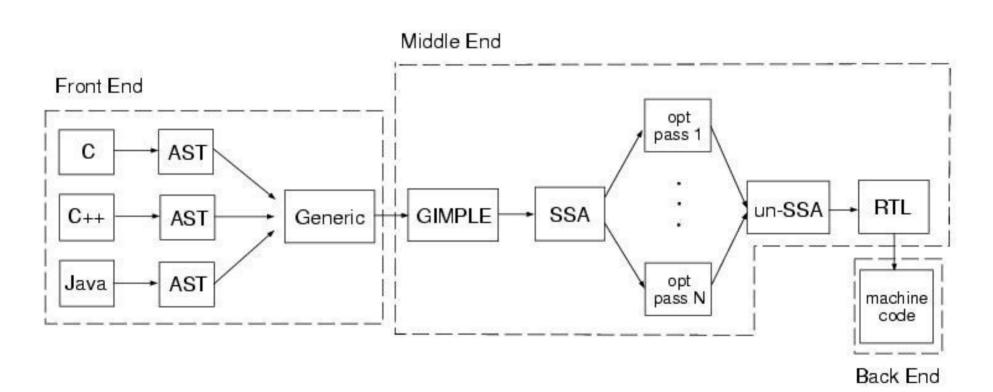
How are ELFs made?

Source → Compiler → Assembler → Linker → Loader

- Compiler: GCC, LLVM/Clang, Intel Compiler, VisualStudio Compiler, and so many more!
- Assembler: Converts assembly language (text) to machine readable code, produces object files
- Linker: Combines one or more objects to create executable, resolves external references and symbols
- Loader: Loads executable and maps it to process memory space, then executes it

The GCC Compiler

A 1 Mio. Foot View!



```
(const int -4 [0xffffffffffffffffff])) [6 op+0 S4 A32])) sqlite3.c:20007 86 {*movsi internal}
    (nil))
(insn 35 34 55 3 (set (reg:DI 0 ax [111])
                            Die bug Output
          (const int 8 [0x8]))) sqlite3.c:20007 217 {*leadi
    (nil)
(insn 37 55 38 3 (set (reg:DI 0 ax [114])
       (symbol ref:DI ("sqlite3Stat") [flags 0x2] <var decl 0x7f71346f4ab0 sqlite3Stat>)) sqlite3.c:20007 85 {*movdi in
    (nil))
(insn 38 37 39 3 (set (reg:DI 1 dx [orig:92 11 ] [92])
       (mem:DI (plus:DI (reg:DI 1 dx [113])
    .... is worth gold, and looks a bit like a "Matrix" screensaver
       when you scroll down fast
              (const int -4 [0xffffffffffffffffff])) [6 op+0 S4 A32])) sqlite3.c:20007 86 {*movsi internal}
    (nil))
(insn 40 39 41 3 (set (reg:DI 0 ax [115])
    (sign extend:DI (reg:SI 0 ax [116]))) sqlite3.c:20007 147 {*extendsidi2_rex64} (nil) fdump-passes
       -fdump-tree-all, -fdump-ipa-all, -fdump-rtl-all
       -fdump-tree-cfg-all
(insn 56 41 43 3 (set regin) 2 CA THE SOMEPLUGIN
```

(mult:DI (reg:DI 0 ax [117])

(nil))

(const int 8 [0x8]))) sqlite3.c:20007 217 {*leadi}

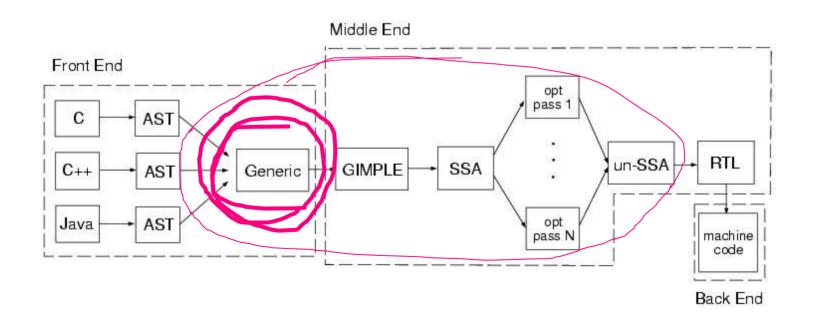
GCC Plugins

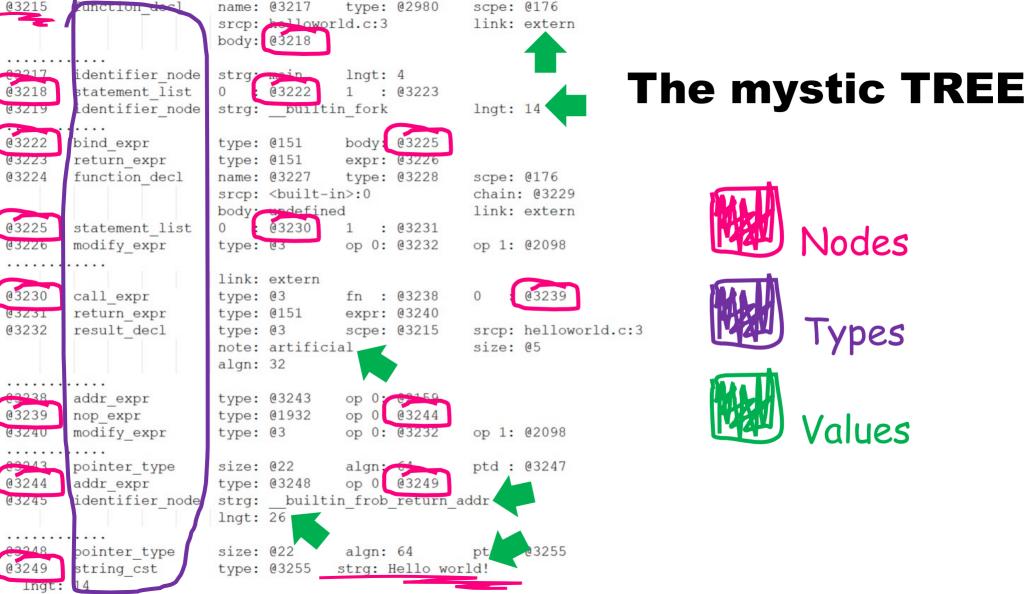
Since GCC 4.5 we can plug passes into the compilation process! Benefits of plugins vs. modifying GCC itself?

- Plugins are shared objects, loaded by GCC as dedicated passes
- Maintained by pass manager
- Dependent on compiler version
- GCC plugin API defined in tree-pass.h

GCC Representations and Data Structures

GENERIC or the mystic TREE







Values

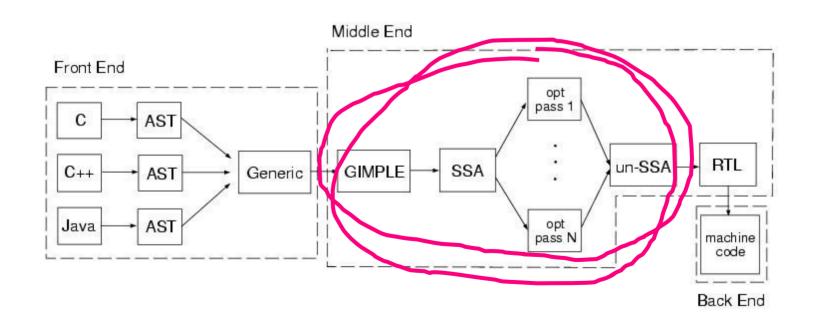
GENERIC

- Language-independent way of representing an entire function as trees
- Interface between parser and optimizers
- Superset of Gimple: imagine a language with a tree structure, similar to LISP
- Defined in gcc/tree.def

Important concepts:

Tree types and DECLs, expressions and statements

GIMPLE – A tree based representation



GIMPLE

The three address code

Target- and language independent optimization

```
# Calculate one solution to the [[quadratic
equation]].
x = (-b + sqrt(b^2 - 4*a*c)) / (2*a)
t1 := b * b
t2 := 4 * a
t3 := t2 * c
t4 := t1 - t3
t5 := sqrt(t4)
t6 := 0 - b
t7 := t5 + t6
t8 := 2 * a
t9 := t7 / t8
x := t9
```

https://en.wikipedia.org/wiki/Three-address_code

```
∃int sub (void) {
         int x = 10;
         int y = 2;
         return x-v;
 5
    ∃int add(void) {
 8
         int a = 10:
 9
         int b = 20;
10
         return a+b:
11
12
13
    ⊟int main (void) {
14
         int d = add();
15
         int f = sub():
16
17
         int q = (d * 100 + 15) - (f * 10 - 50);
18
         return 0:
19
```

20

GIMPLE

26

```
sub ()
      int D.1809;
      int x;
      int y;
      x = 10;
 8
      v = 2;
 9
      D.1809 = x - y;
10
      return D.1809;
11
12
13
14
    add ()
15
16
      int D.1811;
17
      int a;
18
      int b;
19
20
      a = 10;
21
      b = 20;
      D.1811 = a + b;
23
      return D.1811;
24
25
```

```
main ()
28
29
      int D.1813;
30
31
32
        int d;
33
        int f;
34
        int q;
35
36
        d = add();
37
        f = sub();
           = d * 100;
       39
40
        _4 =
            3 + -50;
41
        g = 2 - 4;
42
43
        D.1813 = 0;
44
        return D.1813;
45
46
      D.1813 = 0;
47
      return D.1813;
48
```

GIMPLE in code

Instruction set and language structure much like any high level programming language

GIMPLE_ASSIGN, GIMPLE_CALL, GIMPLE_RETURN, etc.

GIMPLE_PHI, GIMPLE_ASM, etc.

Iterators & statement modifiers

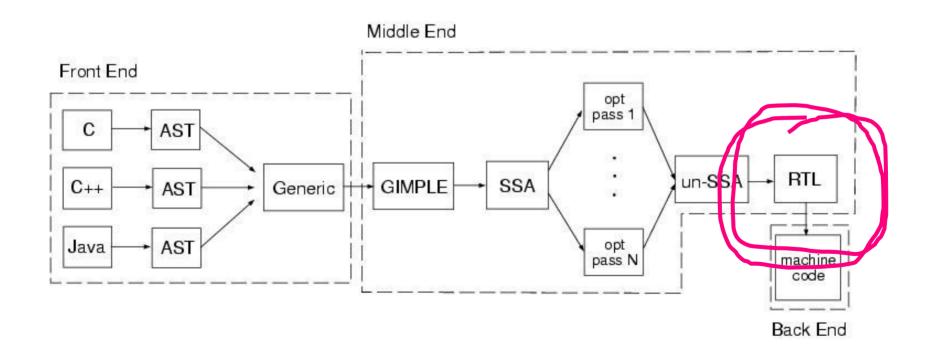
Closely tied to TREE

Go-to tool for CFG and Tree SSA optimizers in GCC middle end

GIMPLE in code

```
// Iterating through basic blocks and Gimple sequences
FOR EACH BB FN (bb, cfun) (
                                                                           Iterator
   for (gsi = gsi start bb(bb); !gsi end p(gsi); gsi next(&gsi)) {
       gimple *statement = gsi stmt(gsi);
       // Picking up on the printf within our helloworld.c
       if (gimple code(statement) == GIMPLE CALL) {
                                                                   Searching CALL statement
           // Getting the first argument of printf
           tree arg = gimple call arg(statement, 0);
           // Building the new string argument
           tree satan = build string(strlen("Hail Satan!!\n")+1, "Hail Satan!!\n");
           tree type = build array type (
               build type variant (char type node, 1, 0),
                                                                                           Building an
               build index type (size int (strlen ("Hail Satan!!\n"))));
           TREE TYPE (satan) = type;
                                                                                           argument
           TREE CONSTANT (satan) = 1;
           TREE READONLY (satan) = 1;
           TREE STATIC (satan) = 1;
           // Replacing the helloworld string argument
           TREE OPERAND (TREE OPERAND ((arg), 0), 0) = satan;
           gimple call set arg(statement, 0, arg);
                                                                    Modification
```

Register Transfer Language



RTL

RTL passes "implement" the machine definition machine definition reflects the processor ABI

target dependent optimization register allocation machine code generation rtl.def, rtl.h, <machine>.md emit-rtl.h

"Assembly language for an abstract machine with infinite registers"

Instructions to be generated are described in an algebraic form that describes what the instruction does

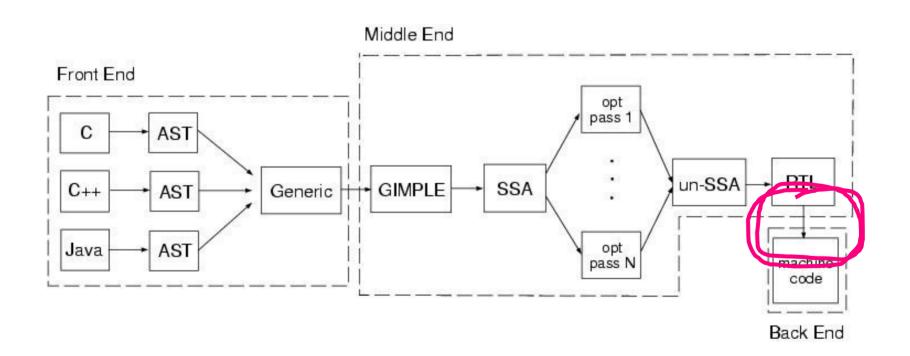
The beauty lies within:)

```
[\ldots]
(insn 5 2 6 2
        (set (reg:DI 5 di)
         (symbol ref/f:DI ("*.LCO") [flags 0x2] <var decl 0x7fd4f1a1ecf0 *.LCO>))
"helloworld.c":4 -1
(nil)
(call insn 6 5 7 2 (set (reg:SI 0 ax)
        (call (mem:QI (symbol ref:DI ("puts") [flags 0x41]
               <function decl 0x7fd4f1974600 builtin puts>) [0 builtin puts S1 A8])
               (const int 0 [0])) "helloworld.c":4 -1
         (nil)
        (expr list:DI (use (reg:DI 5 di))
(nil)))
[\ldots]
```

```
// lea scratchReg, [memLocation + offset]
mySymbol= gen rtx SYMBOL REF(Pmode, memLocation);
SYMBOL REF FLAGS(mySymbol) |= SYMBOL FLAG LOCAL;
leaInstruction = gen rtx SET(scratchReg, plus constant(Pmode, mySymbol, offset));
emit insn before(leaInstruction, positionInsn);
// push variable (64 bit)
decrementStackP = gen_rtx_PRE_DEC(DImode, stack pointer rtx);
topOfStack = gen rtx MEM(DImode, decrementStackP);
                                                             How do we
pushInstruction = gen rtx SET(topOfStack, variable);
emit insn before(pushInstruction, positionInsn);
                                                              generate
                                                             RTL within
// mov scratchReg, sourceReg
                                                                GCC?
movInstruction = gen rtx SET(scratchReg, sourceReg);
```

```
// call <location>
myInternalLabel = gen label rtx();
LABEL NUSES(myInternalLabel)++;
ASM_GENERATE_INTERNAL_LABEL(LNAME, "L", CODE_LABEL_NUMBER(myInternalLabel));
mySymbol = gen_rtx_SYMBOL_REF(Pmode, LNAME);
callInstruction = gen_rtx_CALL(Pmode, gen_rtx_MEM(FUNCTION_MODE,
mySymbol), const0_rtx);
emit call insn before(callInstruction, insn);
[\ldots]
emit label before(myInternalLabel, insnAtLocation);
```

Machine Definitions kmachine.md



Machine Definitions <a href="m

Main part of a gcc backend to be found in gcc/config/<machine>

i386.md

i386.opt

i386-modes.def

i386-protos.h

i386.c and i386.h

IPA – Inter-Procedural Analysis

- IPA passes operate on the call graph and the varpool, inter-procedurally
- IPA_PASS and SIMPLE_IPA_PASS
- IPA LTO: stages partially run at compile time or at link time
- Go-to tools are essentially GIMPLE and GENERIC

```
generate_summary
write_summary
read_summary
execute
write_optimization_summary
read_optimization_summary
function_transform
variable_transform
```

But Security!!1!

Compiler Optimizations

... and what they have to do with security.



How do we get infoleak?

- Convert UAF into OOB read
- Stylesheet that reads attribute name
- Firefox XSLT stores attribute as a Element prt + attribute index
- Free Element and replace with another Element with less attributes

```
void txXPathNodeUtils::getNodeName(const txXPathNode& aNode, nsAString& aName) {
  aNode.Content()
       ->AsElement()
       ->GetAttrNameAt(aNode.mIndex)
       ->GetQualifiedName(aName);
```











OFFENSIVE CON

How do we ge

- Convert UAF in
- Stylesheet that
- Firefox XSLT st
- Free Element a

```
void txXPathNodeUti
...
aNode.Content()
    ->AsElement(
    ->GetAttrNam
    ->GetQualifi
}
```

38:55 / 47:41

OFFENSIVE CON

● □ □

This shouldn't work...

```
const nsAttrName* AttrArray::GetSafeAttrNameAt(uint32_t aPos) const {
  if (aPos >= AttrCount()) {
    return nullptr;
  }
  return &mImpl->mBuffer[aPos].mName;
}
```





How do we ge

- Convert UAF in
- Stylesheet that
- Firefox XSLT st
- Free Element a

```
void txXPathNodeUti
...
aNode.Content()
    ->AsElement(
    ->GetAttrNam
    ->GetQualifi
```

38:55 / 47:41

This shouldn't wo

```
const nsAttrName* At
  if (aPos >= AttrCo
    return nullptr;
  }
  return &mImpl->mBu
}
```

39:54 / 47:41

OFFENSIVE () CO

...and yet it does!

OFFENSIVE() CO

```
const nsAttrName* AttrArray::GetSafeAttrNameAt(uint32_t aPos) const {
   if (aPos >= AttrCount()) {
     return nullptr;
   }
   return &mImpl->mBuffer[aPos].mName;
}
```

- The bounds check gets optimized away. Why?
 - Undefined behavior to the rescue!

Undefined Behavior

- There is a bounds check in the source, which got optimized away in the binary
- Ivan Fratric: "Why on earth would the compiler remove a bounds check?"
 - Bounds check returns nullptr
 - Calling function doesn't check for nullptr, so if check fails that results in null pointer dereferentiation
 - Compiler knows that null pointer deref is undefined behavior; undefined behavior greenlights compiler to take remediative action, and compiler decided to remove the check entirely

Are there other such unicorns?

- OpenSSL "Memset Removal" Bug (CVE-2008-1196)
 - Cryptographic material wasn't properly removed from memory after use since compiler thought operation was obsolete since the data wasn't used after memset
- Clang/LLVM Stack Protector Misoptimization (CVE-2020-10771)
 - Clang removed stack protectors under certain conditions, eg. with tail calls
- Mozilla's ARM64 Ion compiler in Firefox (CVE-2023-29548)
 - Wrong lowering (compiling from higher-level IR to lower level) of unsigned division when divisor is negation of power of two. The optimization incorrectly used absolute value logic
- Mbed TLS (versions prior to 3.6.4) (CVE-2025-52496)
 - Compiler removed security critical code due to a race condition in AESNI detection when certain compiler optimizations are applied, attackers can extract cryptographic keys or forge messages

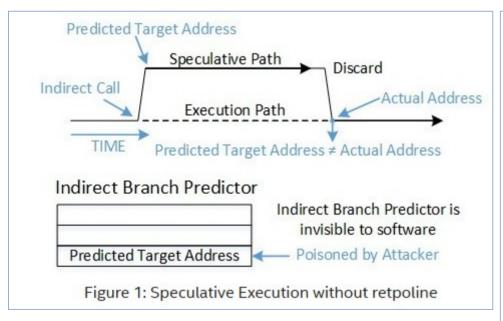
Compiler Mitigations

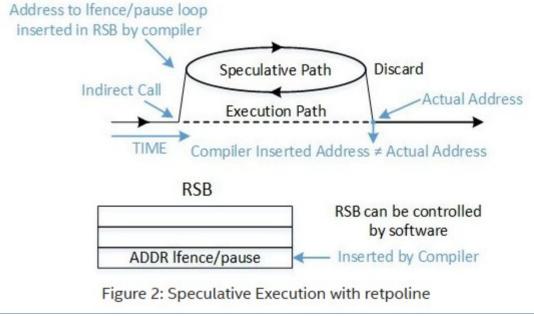
- Stack Protectors
 - Canaries, Safe Stack, Shadow Stack, etc.
- Control Flow Integrity (CFI)
- CPU side channel mitigations
- Address Space Layout Randomization support (ASLR) through position independent code (PIC/PIE)
- And many others...

Spectre v2: Branch Target Injection

- Abuses indirect branch prediction, speculative execution and cache timing side-channels
- Tricks the CPU into speculatively executing memory it wouldn't have executed otherwise, by poisoning indirect branch prediction
- If speculative execution leaves state behind in cache that state can be inferred using cache inference attacks
- This allows attackers to read privileged memory

Spectre v2: Retpoline





Retpoline in GCC

Memory thunk where the function address is at the top of the stack:

```
x86 indirect thunk:
        call L2
L1:
        pause
       1fence
       jmp L1
12:
       lea 8(%rsp), %rsp|lea 4(%esp), %esp
        ret
Indirect jmp via memory, "jmp mem", is converted to
        push memory
        jmp x86_indirect_thunk
Indirect call via memory, "call mem", is converted to
        jmp L2
L1:
       push [mem]
        jmp x86 indirect thunk
L2:
        call L1
```

Register thunk where the function address is in a register, reg:

```
x86 indirect thunk reg:
       call L2
L1:
       pause
       lfence
               L1
L2:
               %reg, (%rsp) | movl %reg, (%esp)
       mova
       ret
where reg is one of (r|e)ax, (r|e)dx, (r|e)cx, (r|e)bx, (r|e)si, (r|e)di,
(r|e)bp, r8, r9, r10, r11, r12, r13, r14 and r15.
Indirect jmp via register, "jmp reg", is converted to
       jmp x86 indirect thunk reg
Indirect call via register, "call reg", is converted to
       call x86 indirect thunk reg
```

Software Protections

- Source level (Developer)
- IR-level (Compiler)
- Machine code / post link stage (aka. Runtime packers)

Compiler Based Protections

- Control flow flattening
- Basic block splitting / reordering
- Bogus control flow / instructions
 - Opaque branches, bloat code, etc.
- String encryption, constant encryption
- Replace simple arithmetic with complex operations
- Breakpoint / hook / debugger detections
- Function inlining/outlining
- Instruction substitution
- SO MUCH MORE

Why Obfuscation Matters

```
.text:00401047
                                  1nc 481847 ·
.text:00401047
                                                                            - CODE XREF: .text:loc 4010471;
.text:00401047 OF 84 FF FF FF(FF
                                                   jz
                                                           near ptr loc 401047+5
.text:0040104D 00 68 43
                                                           [eax+43h], ch
                                                   add
.text:00401050 22 15 90 58 31 05
                                                           a1, a5:5375890n
                                                   and
                                                           [eax+0], dh
.text:00401056 00 30
                                                   add
.text:00401058 40
                                                   inc
                                                           eax
.text:00401059 00 89 00 00 00 10
                                                           [ecx+100000000h], bh
                                                   add
.text:00401045 3B CO
                                                            eax, eax
                                                    CMP
.text:00401045
.text:00401047 OF
                                                    db 0Fh
.text:00401048 84 FF FF FF
                                                    dd 0FFFFFF84h
.text:0040104C
.text:0040104C FF 00
                                                            dword ptr [eax]
                                                    inc
.text:0040104E 6x 10 22 4E 00
                                                            0.64 E 22 h 2 h
                                                    nuch
.text:00401053 58
                                                    pop
                                                            eax
.text:00401054 31 05 00 30 40 00
                                                            dword 403000, eax
                                                    xor
.text:0040105A B9 00 00 00 10
                                                            ecx, 10000000h
                                                    mov
```

Obfuscation Passes **Property** For Developers The items in this section describe the different obfuscation passes available in O-MVLL. The

documentation of the passes provides information about when and how to use the obfuscation as well as

 Dynamic: the pass aims at protecting against dynamic code analysis (hooking, instrumentation, debugging, ...) Control-Flow Flattening Resilience

This property tries to give a level of strength against a fully automated attack. If such an attack exists, it is ✓ Indirect Branch mentioned in the References section of the pass.

片 Indirect Call

12 Function Outline

A Strings Encoding

Obfuscation Overhead

Objective-C Cleaner The kind of overhead introduced in the program when using the given obfuscation pass. The values can be: Opaque Constants

→ Opaque Fields Access

 Code size: the assembly code is larger. · Data size: the raw data of the binary is larger. Memory size: the size of the binary in memory is larger 1.

Thank you

Marion Marschalek pinkflawd@gmail.com



Resources

https://code.woboq.org/gcc/gcc/

https://gcc.gnu.org/onlinedocs/gccint/index.html

https://github.com/enferex/sataniccanary/

https://github.com/ephox-gcc-plugins

https://medium.com/@prathamesh1615/adding-peephole-optimization-to-gcc-89c329dd27b3

https://www.airs.com/dnovillo/200711-GCC-Internals/200711-GCC-Internals-7-passes.pdf

https://www.mitre.org/sites/default/files/publications/supply-chain-attack-framework-14-0228.pdf

https://lwn.net/Articles/457543/

https://www.cse.iitb.ac.in/grc/slides/cgotut-gcc/topic8-retarg-mode.pdf

https://www.cse.iitb.ac.in/~uday/courses/cs715-09/gcc-rtl.pdf

https://en.wikibooks.org/wiki/GNU C Compiler Internals/GNU C Compiler Architecture

https://codesynthesis.com/~boris/blog/2010/05/03/parsing-cxx-with-gcc-plugin-part-1/

https://kristerw.blogspot.com/2017/08/writing-gcc-backend 4.html

https://www.usenix.org/sites/default/files/conference/protected-files/kemerlis usenixsecurity12 slides.pdf

ftp://gcc.gnu.org/pub/gcc/summit/2003/GENERIC%20and%20GIMPLE.pdf

https://pdfs.semanticscholar.org/cafc/c15a1602c5a8090606333b3bdb42e9e80654.pdf

https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/technical-documentation/speculative-execution-side-channel-mitigations.html

https://www.youtube.com/watch?v=U1kc7fcF5Ao