

Vulnerability Modeling with Binary Ninja

Pacific Hackers 2018

Who I am



Josh Watson (@josh_watson)

- Senior Security Engineer for Trail of Bits
- Previously a cog in the Military Industrial Complex
- Before that, other stuff
- Prominent member of the Binary Ninja community

Agenda



1. Introduction
2. Case study: Heartbleed
3. Automated bug hunting without source
4. Results
5. Conclusion
6. Q&A

Static analysis is hard.

Static analysis is harder when you don't have source.

**Hacking like it's 2014:
let's find Heartbleed!**

**TRAIL
OF
BITS**



tls1_process_heartbeat

```
hbtype = *p++;  
n2s(p, payload);  
pl = p;  
/* Skip some stuff... */  
if (hbtype == TLS1_HB_REQUEST)  
{  
    unsigned char *buffer, *bp;  
    int r;  
  
    /* Allocate memory for the response, size is 1 bytes  
     * message type, plus 2 bytes payload length, plus  
     * payload, plus padding  
     */  
    buffer = OPENSSL_malloc(1 + 2 + payload + padding);
```


Heartbleed

tls1_process_heartbeat

```
hbtype = *p++;
n2s(p, payload);
pl = p;
/* Skip some stuff... */
if (hbtype == TLS1_HB_REQUEST)
{
    unsigned char *buffer, *bp;
    int r;

    /* Allocate memory for the response, size is 1 bytes
     * message type, plus 2 bytes payload length, plus
     * payload, plus padding
     */
    buffer = OPENSSL_malloc(1 + 2 + payload + padding);
```

Heartbleed



tls1_process_heartbeat

```
bp = buffer;

/* Enter response type, length and copy payload */
*bp++ = TLS1_HB_RESPONSE;
s2n(payload, bp);
memcpy(bp, pl, payload);
bp += payload;
/* Random padding */
RAND_pseudo_bytes(bp, padding);

r = ssl3_write_bytes(s, TLS1_RT_HEARTBEAT, buffer, 3 + payload + padding);
```

Heartbleed



tls1_process_heartbeat

```
bp = buffer;

/* Enter response type, length and copy payload */
*bp++ = TLS1_HB_RESPONSE;
s2n(payload, bp);
memcpy(bp, pl, payload);
bp += payload;
/* Random padding */
RAND_pseudo_bytes(bp, padding);

r = ssl3_write_bytes(s, TLS1_RT_HEARTBEAT, buffer, 3 + payload + padding);
```

Heartbleed



tls1_process_heartbeat

```
bp = buffer;

/* Enter response type, length and copy payload */
*bp++ = TLS1_HB_RESPONSE;
s2n(payload, bp);
memcpy(bp, pl, payload);
bp += payload;
/* Random padding */
RAND_pseudo_bytes(bp, padding);

r = ssl3_write_bytes(s, TLS1_RT_HEARTBEAT, buffer, 3 + payload + padding);
```

How do we automate analysis of this?

“Let’s be clear: it is trivial to create a static analyzer that runs fast and flags heartbleed. I can accomplish this, for example, by flagging a taint error in every line of code that is analyzed. The task that is truly difficult is to create a static analysis tool that is performant and that has a high signal to noise ratio for a broad range of analyzed programs.”

—John Regehr

A New Development for Coverity and Heartbleed

TRAIL
OF
BITS

Previous work modeling Heartbleed

On detecting Heartbleed with static analysis

Byte-swapping is probably untrusted data that should be tainted.

Look for:

1. byte-swapping
2. combining smaller integers into larger ones
3. byte-swapped values being used as array indices or size parameters for memcpy

```

1. byte_swapping: Performing a byte swapping operation on p implies that it came from an external source, and is therefore tainted.
2. var_assign_var: Assigning: payload = ((unsigned int)p[0] << 8) | (unsigned int)p[1]. Both are now tainted.
2446     n2s(p, payload);
2447     p1 = p;
2448
3. Condition s->msg_callback, taking true branch
2449     if (s->msg_callback)
2450         s->msg_callback(0, s->version, TLS1_RT_HEARTBEAT,
2451             &s->s3->rrec.data[0], s->s3->rrec.length,
2452             s, s->msg_callback_arg);
2453
4. Condition hbtype == 1, taking true branch
2454     if (hbtype == TLS1_HB_REQUEST)
2455     {
2456         unsigned char *buffer, *bp;
2457         int r;
2458
2459         /* Allocate memory for the response, size is 1 bytes
2460          * message type, plus 2 bytes payload length, plus
2461          * payload, plus padding
2462          */
2463         buffer = OPENSSL_malloc(1 + 2 + payload + padding);
2464         bp = buffer;
2465
2466         /* Enter response type, length and copy payload */
2467         *bp++ = TLS1_HB_RESPONSE;
2468         s2n(payload, bp);
2469
2470     CID 1201699 (#1 of 1): Untrusted value as argument (TAINTED_SCALAR)
2471     5. tainted_data: Passing tainted variable payload to a tainted sink.
2472     memcpy(bp, p1, payload);

```

Previous work modeling Heartbleed

Using Static Analysis and Clang To Find Heartbleed

n2s calls are probably untrusted data that should be tainted.

Look for:

1. Results of n2s calls
2. No constraints on the results
3. The results being used as a size parameter in memcpy calls

Modifications to source are needed to facilitate this.

```

2561      /* Read type and payload length first */
2562      hbtype = *p++;
2563      n2s(p, payload);
2564      p += 2;
2565      pl = p;
2566
2567      if (s->msg_callback)
2568          s->msg_callback(0, s->version, TLS1_RT_HEARTBEAT,
2569                          &s->s3->rrec.data[0], s->s3->rrec.length,
2570                          s, s->msg_callback_arg);
2571
2572      if (hbtype == TLS1_HB_REQUEST)
2573      {
2574          unsigned char *buffer, *bp;
2575          int r;
2576
2577          /* Allocate memory for the response, size is 1 bytes
2578             * message type, plus 2 bytes payload length, plus
2579             * payload, plus padding
2580             */
2581          buffer = OPENSSL_malloc(1 + 2 + payload + padding);
2582          bp = buffer;
2583
2584          /* Enter response type, length and copy payload */
2585          *bp++ = TLS1_HB_RESPONSE;
2586          s2n(payload, bp);
2587          memcpy(bp, pl, payload);

```

1 Taking false branch →

2 ← Assuming 'hbtype' is equal to 1 →

3 ← Taking true branch →

4 ← Tainted, unconstrained value used in memcpy size

What do those tools have in common?

**TRAIL
OF
BITS**

Automated bug hunting without source

TRAIL
OF
BITS

Building our toolkit

TRAIL
OF
BITS

Constraint Solving



$$x + y = 8$$

$$2x + 3 = 7$$

Constraint Solving



$$x + y = 8$$

$$2x + 3 = 7$$

$$x = ?$$

$$y = ?$$

Constraint Solving



$$x + y = 8$$

$$2x + 3 = 7$$

$$x = 2$$

$$y = 6$$

Constraint Solving with Z3

```
>>> from z3 import *
>>> x = Int('x')
>>> y = Int('y')
>>> s = Solver()
>>> s.add(x + y == 8)
>>> s.add(2*x + 3 == 7)
>>> s.check()
sat
>>> s.model()
[x = 2, y = 6]
```

Constraint Solving with Z3

```
lea eax, [ebx+8]
cmp eax, 0x20
jle allocate
int3
allocate:
push eax
call malloc
ret
```


sub_0:

0 @ 00000000 `eax = ebx + 8`

1 @ 00000006 `if (eax <= 0x20) then 2 @ 0x9 else 5 @ 0x8`

2 @ 00000009 `push(eax)`
3 @ 0000000a `call(malloc)`
4 @ 0000000f `<return> jump(pop)`

5 @ 00000008 `breakpoint`

Constraint Solving with Z3



`eax = ebx + 8`

`ebx > 0x20`

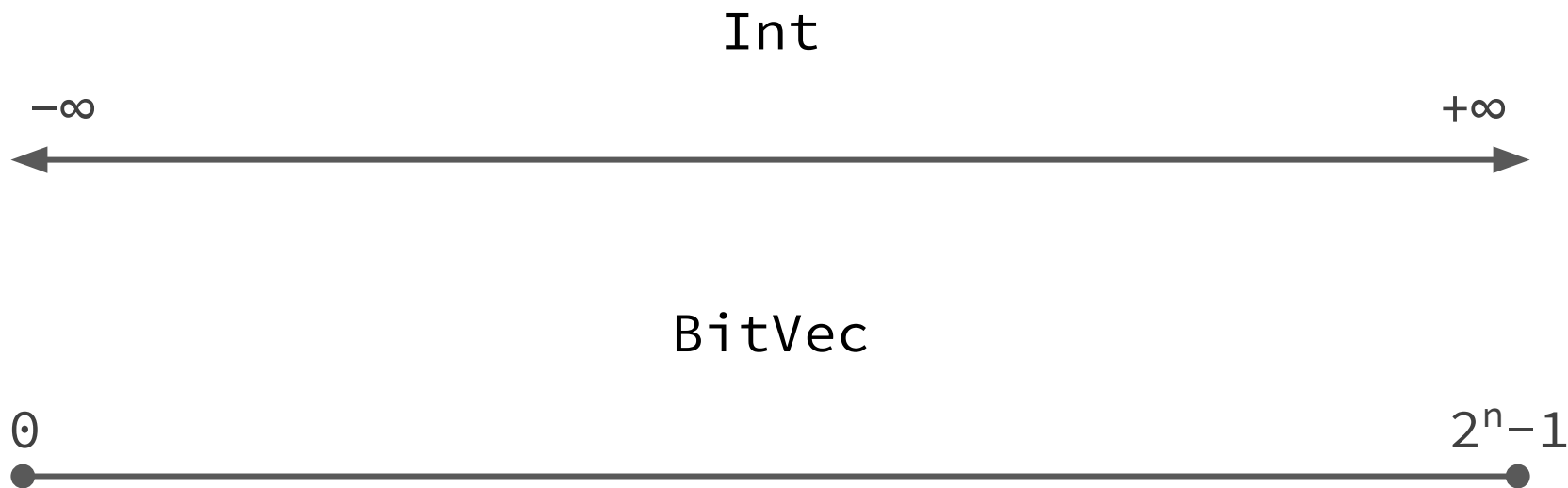
`eax <= 0x20`

Constraint Solving with Z3

```
>>> eax = Int('eax')
>>> ebx = Int('ebx')
>>> s = Solver()
>>> s.add(eax == ebx + 8)
>>> s.add(ebx > 0x20)
>>> s.add(eax <= 0x20)
>>> s.check()
unsat
```

Constraint Solving with Z3

What went wrong?



Constraint Solving with Z3

```
>>> eax = Int('eax')
>>> ebx = Int('ebx')
>>> s = Solver()
>>> s.add(eax == ebx + 8)
>>> s.add(ebx > 0x20)
>>> s.add(eax <= 0x20)
>>> s.check()
unsat
```

Constraint Solving with Z3

```
>>> eax = BitVec('eax', 32)
>>> ebx = BitVec('ebx', 32)
>>> s = Solver()
>>> s.add(eax == ebx + 8)
>>> s.add(ebx > 0x20)
>>> s.add(eax <= 0x20)
>>> s.check()
sat
```

Constraint Solving with Z3

```
>>> s.model()
[ebx = 2147483640, eax = 2147483648]
>>> hex(2147483640)
'0x7fffffff8'
>>> hex(2147483648)
'0x800000000'
```

Constraint Solving with Z3



How do we collect semantics and translate those into a set of constraints?

- ❑ BitVec trivially works for registers
- ❑ Memory accesses need Z3's Array sort
- ❑ Stack variables reside in memory, tracking pushes and pops is complicated
- ❑ Can we treat stack variables the same as registers?



```

sub_0:
00000000 push    ebp
00000001 mov     ebp, esp
00000003 sub     esp, 0x8 {var_c}
00000006 mov     dword [ebp-0x4 {var_8}], 0x0
0000000d mov     dword [ebp-0x8 {var_c}], 0x10
00000014 mov     eax, dword [ebp+0x8 {arg1}]
00000017 add     eax, dword [ebp-0x8 {var_c}]
0000001a mov     esp, ebp
0000001c pop     ebp
0000001d retn

```

LLIL

```

sub_0:
0 @ 00000000 push(ebp)
1 @ 00000001 ebp = esp {__saved_ebp}
2 @ 00000003 esp = esp - 8
3 @ 00000006 [ebp - 4 {var_8}].d = 0
4 @ 0000000d [ebp - 8 {var_c}].d = 0x10
5 @ 00000014 eax = [ebp + 8 {arg1}].d
6 @ 00000017 eax = eax + [ebp - 8 {var_c}].d
7 @ 0000001a esp = ebp
8 @ 0000001c ebp = pop
9 @ 0000001d <return> jump(pop)

```

MLIL

```

sub_0:
0 @ 00000014 int32_t eax = arg1
1 @ 00000017 int32_t eax_1 = eax + 0x10
2 @ 0000001d return eax_1

```

```

int64_t __saved_rbp {Frame offset -8}
void* const __return_addr {Frame offset 0}
int64_t result {Register rax}
struct LIST_ENTRY* Next {Register rax}
struct LIST_ENTRY* Next {Register rax}
struct LIST_ENTRY* New {Register rsi}
struct LIST_ENTRY* Prev {Register rdi}

```

_Insert:

```

0 @ 10000ef4 int64_t result = 0
1 @ 10000ef9 if (Prev == 0) then 2 @ 0x10000f2c else 3 @ 0x10000efe

```

```

3 @ 10000efe if (New == 0) then 2 @ 0x10000f2c else 4 @ 0x10000f00

```

```

4 @ 10000f00 struct LIST_ENTRY* Next = Prev->Flink
5 @ 10000f06 if (Next == 0) then 6 @ 0x10000f18 else 9 @ 0x10000f08

```

```

6 @ 10000f18 New->Flink = 0
7 @ 10000f1f New->Blink = Prev
8 @ 10000f1f goto 14 @ 0x10000f23

```

```

9 @ 10000f08 New->Flink = Next
10 @ 10000f0b New->Blink = Prev
11 @ 10000f0f struct LIST_ENTRY* Next = Prev->Flink
12 @ 10000f12 Next->Blink = New
13 @ 10000f16 goto 14 @ 0x10000f23

```

```

14 @ 10000f23 Prev->Flink = New
15 @ 10000f26 int64_t result = 1
16 @ 10000f26 goto 2 @ 0x10000f2c

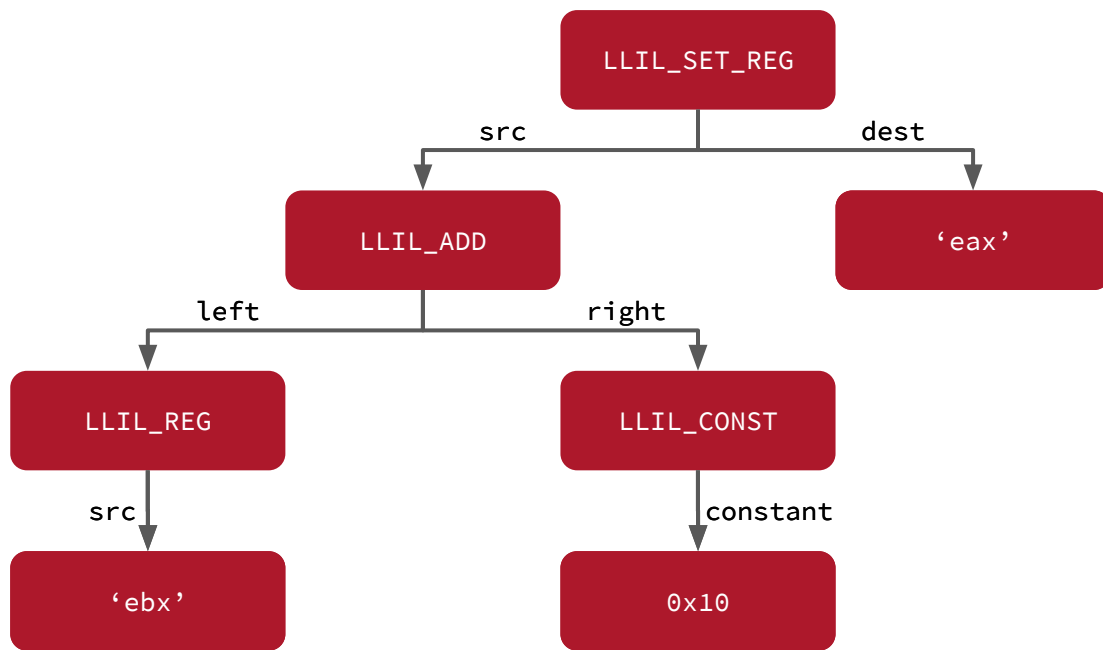
```

```

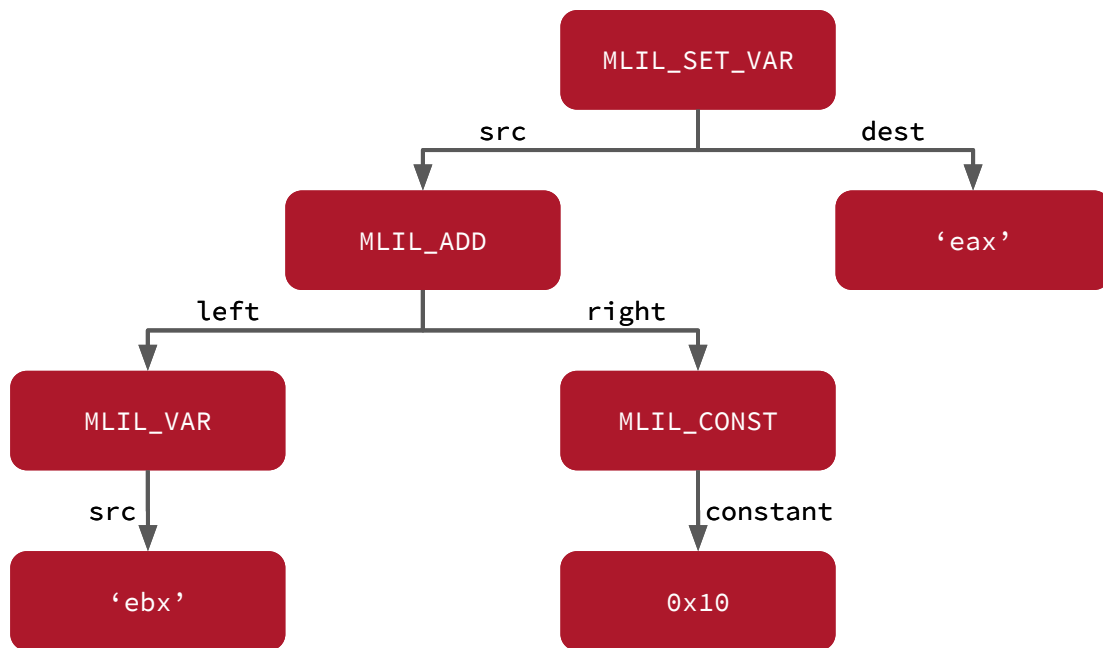
2 @ 10000f2c return result

```

```
lea eax, [ebx + 0x10]
```



`lea eax, [ebx + 0x10]`

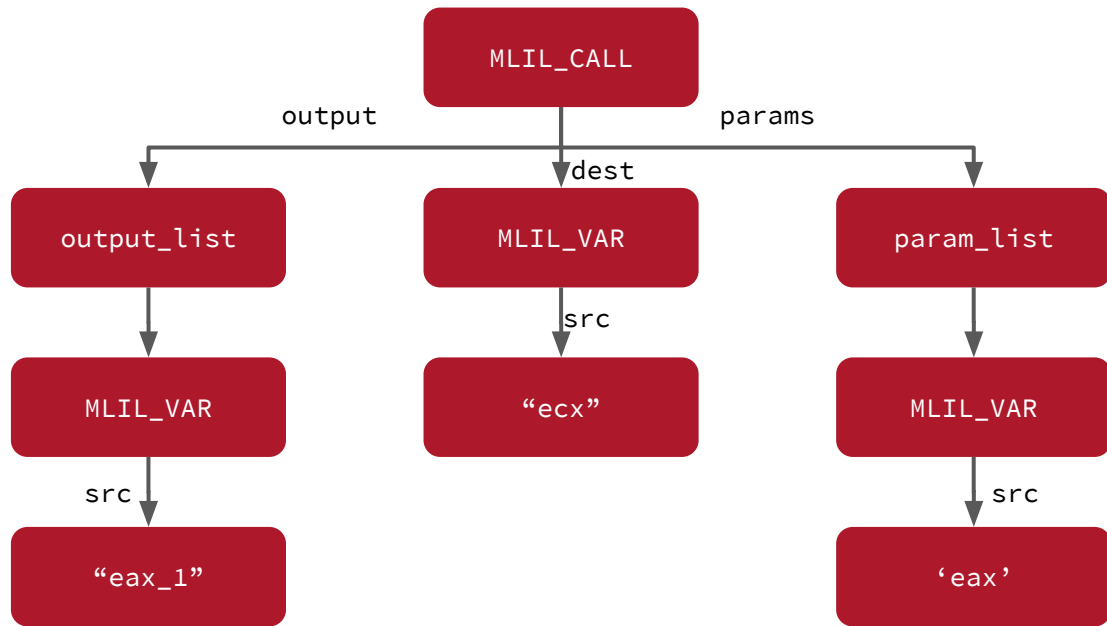


`lea eax, [ebx + 0x10]`

```
push eax; call ecx
```



`push eax; call ecx`



push eax; call ecx

But...

TRAIL
OF
BITS

What's wrong with modeling this?



```
mov  eax, ebx  
lea  eax, [ecx+eax*4]
```

What's wrong with modeling this?

`eax = ebx`

`eax = ecx + eax << 2`

What's wrong with modeling this?

Constraints are purely expressing mathematical truths about variables in a system of equations and have **no** temporal element at all.

Single Static Assignment

Single Static Assignment (SSA)

SSA form is a representation of a program in which every variable is defined once and only once.

If the variable is assigned a new value, a new “version” of that variable is defined instead.

Single Static Assignment

Original form

```
a = 1
b = 2
a = a + b
```

SSA form

```
a1 = 1
b1 = 2
a2 = a1 + b1
```


Single Static Assignment

Original form

```
def f(a):  
    if a > 20:  
        a = a * 2  
    else:  
        a = a + 5  
    return a
```

SSA form

```
def f(a0):  
    if a0 > 20:  
        a1 = a0 * 2  
    else:  
        a2 = a0 + 5  
    a3 =  $\Phi(a_1, a_2)$   
    return a3
```

SSA makes it easy to explicitly track all definitions and uses of a variable throughout the lifetime of the program.

```
sub_0:  
0 @ 00000002 rsi = arg1  
1 @ 00000005 rdi = arg2  
2 @ 0000000e if (arg1 u> 0) then 3 @ 0x12 else 5 @ 0x10
```

```
3 @ 00000012 rax = arg2 << 2  
4 @ 00000012 goto 6 @ 0x1c
```

```
5 @ 00000010 goto 6 @ 0x1c
```

```
6 @ 0000001c return
```

```
sub_0:  
0 @ 00000002 rsi#1 = arg1#0  
1 @ 00000005 rdi#1 = arg2#0  
2 @ 0000000e if (arg1#0 u> 0) then 3 @ 0x12 else 5 @ 0x10
```

```
3 @ 00000012 rax#1 = arg2#0 << 2  
4 @ 00000012 goto 6 @ 0x1c
```

```
5 @ 00000010 goto 6 @ 0x1c
```

```
6 @ 0000001c rax#2 =  $\phi$ (rax#0, rax#1)  
7 @ 0000001c return
```

MLIL SSA form

Medium Level IL

MLIL_SET_VAR

MLIL_VAR

MLIL_CALL

MLIL_LOAD

Variable

Medium Level IL SSA

MLIL_SET_VAR_SSA

MLIL_VAR_SSA

MLIL_CALL_SSA

MLIL_LOAD_SSA

SSAVariable

Putting it all together

TRAIL
OF
BITS

Combining MLIL SSA with Z3



`eax = ebx`

`eax = ecx + (eax << 2)`

Combining MLIL SSA with Z3



`eax#1 = ebx#0`

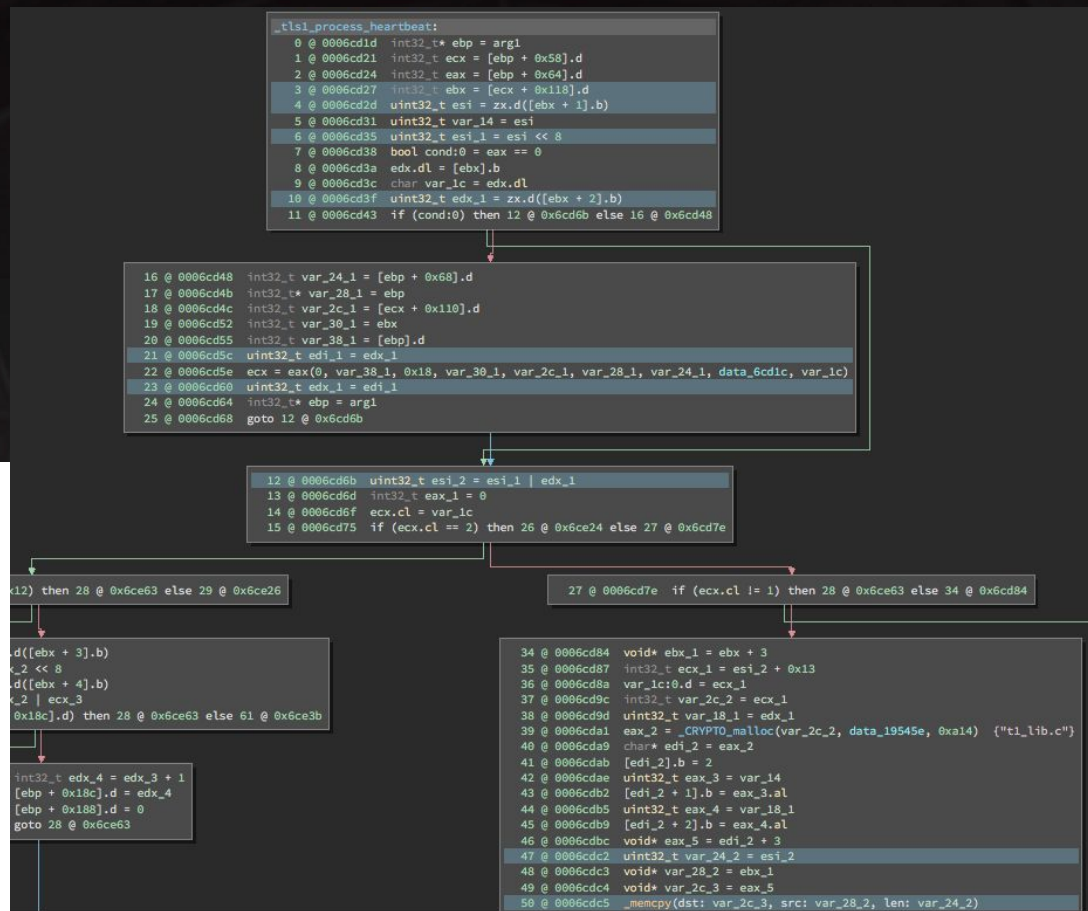
`eax#2 = ecx#0 + (eax#1 << 2)`

Combining MLIL SSA with Z3

```
eax_1 = BitVec('eax#1', 32)
ebx_0 = BitVec('ebx#0', 32)
ecx_0 = BitVec('ecx#0', 32)
eax_2 = BitVec('eax#2', 32)
s = Solver()
s.add(
    eax_1 == ebx_0,
    eax_2 == ecx_0 + (eax_1 << 2)
)
```


Writing the Script

TRAIL
OF
BITS



Writing the script



1. Find our “sinks”
2. Eliminate sinks that are obviously not vulnerable
3. Trace the variables the size depends on (backwards slice)
4. Identify variables that might be part of a byte swap
5. Identify additional constraints on the size parameter
6. Solve the model
7. Find bugs

Writing the script

Step 1: Finding our “sinks”

```
memcpy_refs = [  
    (ref.function, ref.address)  
    for ref in bv.get_code_refs(bv.symbols['_memcpy'].address)  
]  
  
dangerous_calls = []  
  
for function, addr in memcpy_refs:  
    call_instr = function.get_low_level_il_at(addr).medium_level_il  
    if check_memcpy(call_instr.ssa_form):  
        dangerous_calls.append((addr, call_instr.address))
```

Writing the script

Step 2: Eliminate sinks that we know aren't vulnerable

```
def check_memcpy(memcpy_call):  
    size_param = memcpy_call.params[2]  
  
    if size_param.operation != MediumLevelILOperation.MLIL_VAR_SSA:  
        return False  
  
    possible_sizes = size_param.possible_values  
  
    if possible_sizes.type != RegisterValueType.UndeterminedValue:  
        return False  
  
    model = ByteSwapModeler(size_param, bv.address_size)  
  
    return model.is_byte_swap()
```

Step 3: Trace the variables the size depends on

```
var_def = self.function.get_ssa_var_definition(self.var.src)

# Visit statements that our variable directly depends on
self.to_visit.append(var_def)

while self.to_visit:
    idx = self.to_visit.pop()
    if idx is not None:
        self.visit(self.function[idx])
```

Writing the script

Step 3: Trace the variables the size depends on

```
def visit_MLIL_ADD(self, expr):  
    left = self.visit(expr.left)  
    right = self.visit(expr.right)  
  
    if None not in (left, right):  
        return left + right
```

Writing the script

Step 4: Identify variables that might be part of a byte swap

```
def visit_MLIL_VAR_SSA(self, expr):  
    if expr.src not in self.visited:  
        var_def = expr.function.get_ssa_var_definition(expr.src)  
  
        if var_def is not None:  
            self.to_visit.append(var_def)  
  
    src = create_BitVec(expr.src, expr.size)  
  
    value_range = identify_byte(expr, self.function)  
    if value_range is not None:  
        self.solver.add(Or(src == 0, And(src = value_range.step)))  
        self.byte_vars.add(expr.src)  
  
    return src
```


Writing the script

Step 4: Identify variables that might be part of a byte swap

```
phi_values = []

for var in expr.src:
    if var not in self.visited:
        var_def = self.function.get_ssa_var_definition(var)
        self.to_visit.append(var_def)

    src = create_BitVec(var, var.var.type.width)

    # ...

    phi_values.append(src)

if phi_values:
    phi_expr = reduce(
        lambda i, j: Or(i, j), [dest == s for s in phi_values]
    )

    self.solver.add(phi_expr)
```

Writing the script

Step 4: Identify variables that might be part of a byte swap

```
# If this value can never be larger than a byte,
# then it must be one of the bytes in our swap.
# Add it to a list to check later.
if src is not None and not isinstance(src, (int, long)):
    value_range = identify_byte(expr.src, self.function)
    if value_range is not None:
        self.solver.add(Or(src == 0, And(src <= value_range.end, src >= value_range.step)))

    self.byte_vars.add(*expr.src.vars_read)

if self.byte_values.get((value_range.step, value_range.end)) is None:
    self.byte_values[
        (value_range.step, value_range.end)
    ] = simplify(Extract(
        int(math.floor(math.log(value_range.end, 2))),
        int(math.floor(math.log(value_range.step, 2))),
        src
    ))
)
```

Step 5: Identify constraints on the size parameter

```
for i, branch in self.var.branch_dependence.iteritems():
    for vr in self.function[i].vars_read:
        if vr in self.byte_vars:
            raise ModelIsConstrained()

    vr_def = self.function.get_ssa_var_definition(vr)
    if vr_def is None:
        continue

    for vr_vr in self.function[vr_def].vars_read:
        if vr_vr in self.byte_vars:
            raise ModelIsConstrained()
```

Writing the script

Step 6: Solve the model

```
self.solver.add(
    Not(
        And(
            var == ZeroExt(
                var.size() - len(ordering)*8,
                Concat(*ordering)
            ),
            reverse_var == ZeroExt(
                reverse_var.size() - reversed_ordering.size(),
                reversed_ordering
            )
        )
    )
)

if self.solver.check() == unsat:
    return True
```

Results

TRAIL
OF
BITS

Results

```
Blog Post — -bash — 80x21
[(venv) air:Blog Post user$ PYTHONPATH=/Applications/Binary\ Ninja.app/Contents/Resources/python python find_heartbleed.py openssl.bndb
openssl.bndb loaded...
Checking 319 memcpy calls
100% ██████████ 319/319 [00:44<00:00, 2.83it/s]
***** POTENTIAL VULNERABILITY *****
memcpy at 0x78104 in _dtls1_process_heartbeat uses a size parameter that potentially comes from an untrusted source!
***** POTENTIAL VULNERABILITY *****
memcpy at 0x6cdc5 in _tls1_process_heartbeat uses a size parameter that potentially comes from an untrusted source!
Analysis complete.
(venv) air:Blog Post user$
```

OpenSSL 1.0.1f

Compiled with `./Configure darwin-i386-cc`

Two functions are found:

`tls1_process_heartbeat` and
`dtls1_process_heartbeat`.

```
Blog Post — -bash — 80x21
[(venv) air:Blog Post user$ PYTHONPATH=/Applications/Binary\ Ninja.app/Contents/Resources/python python find_heartbleed.py ../Repos/openssl-1.0.1g/openssl.bndb
../Repos/openssl-1.0.1g/openssl.bndb loaded...
Checking 318 memcpy calls
100% ██████████ 318/318 [00:42<00:00, 3.12it/s]
Analysis complete.
(venv) air:Blog Post user$
```

OpenSSL 1.0.1g

Compiled with `./Configure darwin-i386-cc`

The vulnerable functions are no longer present.

Conclusion

TRAIL
OF
BITS

https://github.com/trailofbits/binjascripts/tree/master/find_heartbleed

Binary Ninja's MLIL and SSA form make advanced binary analysis easier to implement.

Questions?



Josh Watson

Senior Security Engineer

josh@trailofbits.com

@josh_watson

References/Links

On detecting Heartbleed with static analysis

<https://www.synopsys.com/blogs/software-security/detecting-heartbleed-with-static-analysis/>

Using Static Analysis and Clang To Find Heartbleed

<https://blog.trailofbits.com/2014/04/27/using-static-analysis-and-clang-to-find-heartbleed>

Heartbleed and Static Analysis


<https://blog.regehr.org/archives/1125>

A New Development for Coverity and Heartbleed

<https://blog.regehr.org/archives/1128>

Vulnerability Modeling with Binary Ninja

<https://blog.trailofbits.com/2018/04/04/vulnerability-modeling-with-binary-ninja/>



TRAIL *OF* **BITS**