

CSci 5271 Guest Lecture

Navid Emamdoost
navid@cs.umn.edu



UNIVERSITY OF MINNESOTA
Driven to Discover™

Software-based Fault Isolation



UNIVERSITY OF MINNESOTA
Driven to Discover™

Need for extensibility

- Applications can incorporate independently developed modules
 - Operating System
 - Add new file system
 - Database Management System
 - User-defined data type
 - Browser
 - Multimedia editor



UNIVERSITY OF MINNESOTA



3

Problem with extensions

- Security and Reliability
- Extensions may be
 - Malicious
 - Vulnerable
 - Faulty
- Solution:
 - Isolate from others



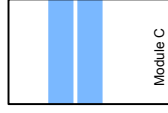
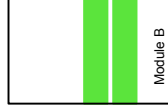
UNIVERSITY OF MINNESOTA



4

Isolation option 1

- Hardware-based isolation
 - Place each module in its own address space
 - Communicate via RPC



RPC

Module A

Module B

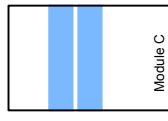
Module C



UNIVERSITY OF MINNESOTA

Isolation option 1

- Hardware-based isolation
 - Place each module in its own address space
 - Communicate via RPC
 - Switch to kernel mode
 - Copy arguments
 - Save/Restore registers
 - Switch address spaces
 - Return to user mode



RPC

Module B

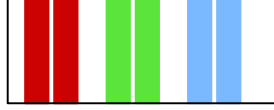
Module C

5

6

Isolation option 2

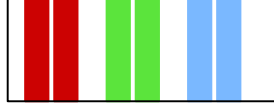
- Software-based isolation
 - All modules in same virtual address
 - Protect them from each other
 - Provide efficient communication



7

Isolation option 2

- Software-based isolation
 - All modules in same virtual address
 - Protect them from each other
 - Provide efficient communication



8

Efficient Software-based Fault Isolation

Robert Wahbe, Steven Lucco, Thomas E. Anderson, Susan L. Graham
SOSP 1993



UNIVERSITY OF MINNESOTA
Driven to Discover™

10

Fault Domain

- Load untrusted extension into its own fault domain
 - Code Segment
 - Data Segment
- Security Policy:
 - No code is executed outside of fault domain
 - No data changed outside of fault domain
 - Some protect load, too



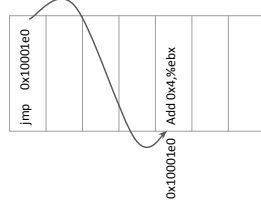
11

Goal

- Protect the rest of an application from a buggy/malicious module on **RISC** architecture
- Separate untrusted code
 - Define a fault domain
 - Prevent the module from jumping or writing outside of it
 - While letting efficient communications

Unsafe Instructions

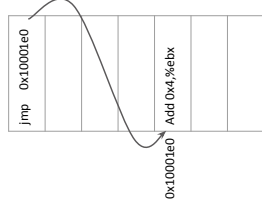
- Jump or store instructions
 - Change Control flow
 - Change data
- Addressing issue
 - jmp 0x10001e0



12

Unsafe Instructions

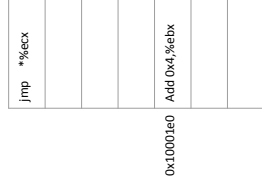
- Jump or store instructions
 - Change Control flow
 - Change data
- Addressing issue
 - jmp 0x10001e0
 - mov %eax, 0x11020028



13

Unsafe Instructions

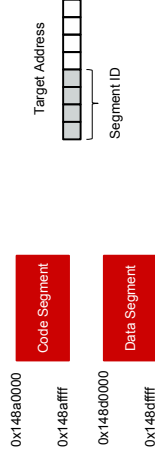
- Jump or store instructions
 - Change Control flow
 - Change data
- Addressing issue
 - jmp 0x10001e0
 - jmp %ecx
 - mov %eax, 0x11020028
 - mov \$0x1b80, (%ecx)



14

Segment ID

- Within a segment
 - Addresses share unique pattern of upper bits



Segment Matching

- Insert checking code before unsafe instruction
 - check segment ID of target address
- Use dedicated registers

dedicated-reg \leftarrow target-address
scratch-reg \leftarrow (dedicated-reg >> shift-reg)
 if scratch-reg == segment-reg:
 jmp/mov using dedicated-reg

Segment Matching

- Needs 4 dedicated registers
- Checking code must be atomic
- Exact location of fault can be detected
- Runtime overhead
 - 4 extra instructions

Address Sandboxing

- Ensure, do not check!
- Before each unsafe instruction
 - Set upper bit of target address to correct segment ID

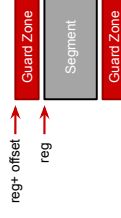
dedicated-reg \leftarrow target-address & and-mask
dedicated-reg \leftarrow dedicated-reg | segment-reg
 jmp/mov using dedicated-reg

Address Sandboxing

- Prevents faults
- Needs 5 dedicated registers
- 2 extra instructions
 - less overhead compared to segment matching

Optimizations

- register-plus-offset mode
 - store value, offset(reg)
 - offset is in the range of -64K to +64K
- mov %esi,0x8(%edx)

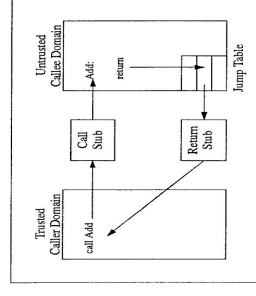


Optimizations

- Stack pointer
 - Just sandbox it when it is set (beginning of a function)
 - Ignore sandboxing for small changes (push, pop)
 - Works because of guard zones

Cross Fault Domain Communication

- Host to an untrusted module
- Untrusted module to host
- Call/Return Stub
 - For each pair of fault domains



22

Cross Fault Domain Call

- Trusted call/return stub
 - copy parameters
 - switch execution stack
 - maintain values of CPU registers
 - no traps or address space switching
 - faster
 - returns via jump table
 - jump targets are immediates
 - a legal address in target fault domain

Implementation

- Change the compiler
 - emit encapsulation code into distrusted code
- At the load time
 - check the integrity of encapsulation code
 - Verifier

24

Verifier

- Responsible for checking encapsulation instructions just before execution start
- Challenge:
 - indirect jump
- Hint:
 - every store/jump uses dedicated registers
- Look for changes in dedicated registers
 - any change means beginning of a check region
 - verify the integrity of check region

25

What about CISC architectures?!

x86

26

Evaluating SFI for a CISC Architecture (PittSField)

Stephen McCamant, Greg Morrisett
USENIX 2005



UNIVERSITY OF MINNESOTA
Driven to Discover™

28

CISC Architectures

- Processor can jump to any byte
- Hard to make hidden instructions safe
- Solution: Instruction Alignment

```

push %esi
mov $0x56,%dh | sbb $0xff,%al | inc %eax | or %al,%dh
movzbl 0x1c(%esi),%edx | incl 0x8(%eax) ...
0f b6      56      1c      ff      40      08      c6
  
```

29

CISC Architectures

- RISC Architecture
 - Fixed length instructions
 - More CPU registers
- Intel IA-32 (aka x86-32)
 - Variable length instructions
 - Less CPU registers
- Classical SFI is not applicable here

28

Alignment

- Divide memory into 16-byte chunks
- No instruction is allowed to cross chunk boundary
- Target of jumps placed at the beginning of chunks
- Call instructions placed at the end of chunk

30

Alignment

- Use *NOP* for padding
- No separation of an unsafe chunk

0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
neg %edi	add \$0x20,%esp	5-byte nop	7-byte nop	7-byte nop	and \$0x10ffff0,%ebx	jmp %ebx	sub %ecx,%eax	9-byte nop	9-byte nop	mov %eax,%eax	5-byte nop	call 0xf50	5-byte nop	call 0xf50	5-byte nop
f10:															
f20:															
f30:															
f40:															
f50:															
f60:															

31

Jumps

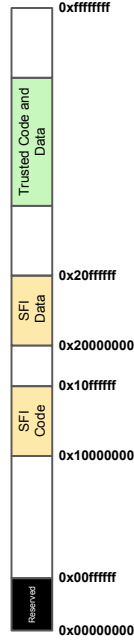
- Chunks are atomic
- Jump destinations are **checked**

0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
neg %edi	add \$0x20,%esp	5-byte nop	7-byte nop	7-byte nop	and \$0x10ffff0,%ebx	jmp %ebx	sub %ecx,%eax	9-byte nop	9-byte nop	mov %eax,%eax	5-byte nop	call 0xf50	5-byte nop	call 0xf50	5-byte nop
f10:															
f20:															
f30:															
f40:															
f50:															
f60:															

32

Optimization: AND-only Sandboxing

- Choose code and data region addresses carefully
- Their ID just has one bit set
- Reduces sandboxing sequence to just one instruction



33

Example

0	1	2	3	4	5	6	7	8	9	a	b	c	d	e	f
neg %edi	add \$0x20,%esp	5-byte nop	7-byte nop	7-byte nop	and \$0x10ffff0,%ebx	jmp %ebx	sub %ecx,%eax	9-byte nop	9-byte nop	mov %eax,%eax	5-byte nop	call 0xf50	5-byte nop	call 0xf50	5-byte nop
f00:															
f10:															
f20:															
f30:															
f40:															
f50:															
f60:															

34

Verification

- Statically check
 - No jump to outside of code region
 - No store to outside of data region
- Before each unsafe jump or store there should be a sandboxing AND
- The sandboxing AND should not be the last instruction in a chunk

35

Performance overhead

- Implemented prototype
 - named PittSField
- Average module overhead: 21%
- But the overall execution can be improved because of faster communications
 - no trap, RPC, etc

Native-client: A Sandbox for Portable, Untrusted x86 Native Code

Bonnet Yee, et al.
IEEE S&P, 2009



UNIVERSITY OF MINNESOTA
Driven to Discover™

38

Google Native Client



UNIVERSITY
OF MINNESOTA

- Browser Plugin (Google Chrome)
 - Allows execution of untrusted C/C++ code in browser
- Browser?! Native Code?!
 - Yes! browsers are new platform for applications
- Gives Browser plugins performance of native code
- Ships by default since Chrome 14

Sandboxing



UNIVERSITY OF MINNESOTA

- Inner Sandbox
 - Code sandboxing
 - Alignment and address sandboxing
 - Check branch target addresses
 - Data Sandboxing
 - segmented addressing mode supported by x86_32
- Outer Sandbox
 - Controls system calls issued by native code
 - Whitelist

39

Inner Sandbox



UNIVERSITY OF MINNESOTA

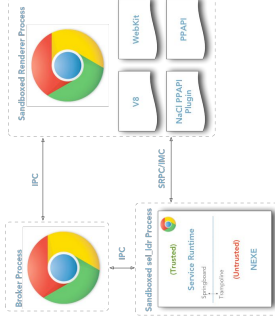
- On x86_32
 - Sandboxing via segmented memory
 - Used to separate trusted from untrusted code/data
 - Modified when switching between trusted/untrusted
 - %cs code
 - %ds data
 - %gs thread local storage
 - %ss %es %fs all set to %ds
- On x86_64
 - mov/branch alignment, guard pages
 - r15 keeps base address of an aligned 4GB range

41

NaCl Architecture



UNIVERSITY OF MINNESOTA



Source: https://media.chromium.orbit.com/view/2284646/NaCl_UIS_12_Rent_Google_NaCl_Sides.pdf

40

Native Client Toolchain



UNIVERSITY OF MINNESOTA

- Modified GCC and GAS
 - To emit sandboxing instructions
- Final executable has ELF file structure (called NEXE)
 - Can be disassembled using standard tools
 - objdump -d

```
nacl.call    %ebx    → and $0xffffffe0,%ebx  
              call    *%ebx  
nacl.ljmp   %ecx    → and $0xffffffe0,%ecx  
              jmp    *%ecx
```

42

Alignment

- Divide memory into 32-byte Bundles
- Target of jumps placed at the beginning of bundles
- No instruction is allowed to cross bundle boundary

```

1000951: 83 48 01          mov     $0x1,%eax
1000954: 88 41 28          mov     %eax,%eax
1000957: 8b 01 04          test   %eax,%eax
100095c: 85 42          jnz     $0x0
1000961: 8b 01 01 00 00    mov     %eax,%eax
1000966: 8b 01 01 00 00    mov     %eax,%eax
1000969: 8d 2c 95 00 00 00 lea    0x0(%eax,%eax),%edi
1000971: 83 09 01          sub    $0x1,%eax
1000974: 83 06 07          and    %eax,%eax
1000977: 8d 06 00 00 00 00 lea    0x0(%eax,%eax),%eax
1000982: 87 06 00 00 00 00 lea    0x0(%eax,%eax,%eax,%eax),%eax
1000987: 8f 06 01 00 00    jmp     $0x0
    
```

43

Alignment

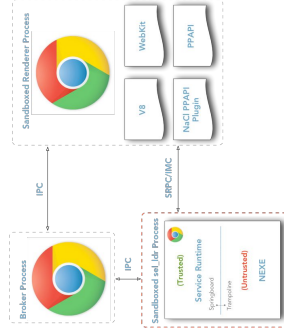
- Divide memory into 32-byte Bundles
- Call instructions placed at the end of bundles
 - For the sake of return address alignment

```

1000700: 55          push   %ebp
1000703: 89 05 18          mov     $0x18,%ebp
1000706: c7 08 2c 2c 56 01 11 movl   $0x11082c2c,%eax
1000709: 8d 04 2c 00 00 00 00 lea    0x0(%eax,%eax),%edi
1000712: c8 00 00 00 00 00 00 scasd  %eax
1000715: 48          call   1004590b
1000720: 89 0c          mov     %ebp,%eax
1000723: 5d          pop    %ebp
1000726: 5b          and   %ebx,%eax
1000729: 99          rep    stc
1000732: 9f 01 e0          and   %eax,%eax
1000735: 9f 01          rep    stc
1000738: 9f 01          rep    stc
1000741: 9f 01          rep    stc
1000744: 9f 01          rep    stc
1000747: 9f 01          rep    stc
    
```

44

Service Runtime



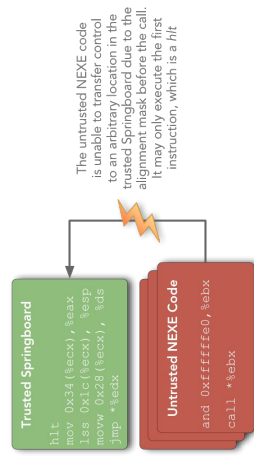
45

Service Runtime

- No existing or new memory allocations may be marked as executable at runtime
 - This guarantees only validated code pages have executable permissions
- NaCl syscall
 - ~ 30 syscalls allowed
 - Basic operations such as open, close, read, write, ioctl, mmap, munmap, stat, _exit and a few others

47

Service Runtime



46

Validator

- Disassembles all NEXE instructions
- Starts at trusted 32 byte aligned entry point
- Exits on any blacklisted instructions
 - Privileged instructions
 - Instructions that modify segment registers
 - ret
 - sysenter

48

CBI NaCl

or
Cross-Bundle Instruction Native Client



UNIVERSITY OF MINNESOTA
Driven to Discover™

49



Padding vs Performance

- Change NaCl padding scheme
 - Pad removal
 - Greedy Algorithm
- Multipass Validator
 - We must guarantee sandboxing policy enforcement
 - Appropriate changes in validator

50



Types of Padding

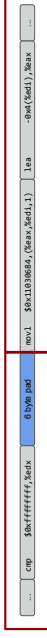
- Indirect jump target
 - Will be placed at the next bundle start
- Call instruction
 - Will be placed at the end of the bundle
- Cross bundle instruction
 - Will be pushed to the start of next bundle

51



Types of Padding

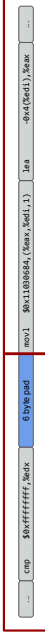
- Indirect jump target
 - Will be placed at the next bundle start
- Call instruction
 - Will be placed at the end of the bundle
- Cross bundle instruction
 - Will be pushed to the start of next bundle
 - Conservative



52



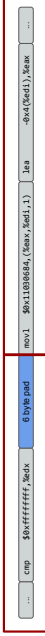
Pad Removal



53

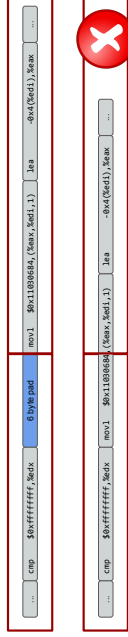


Pad Removal



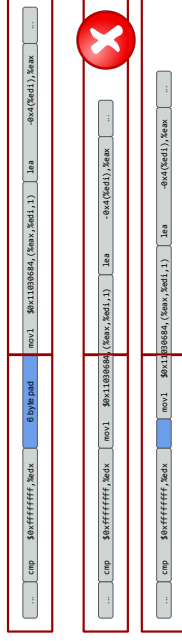
54

Pad Removal



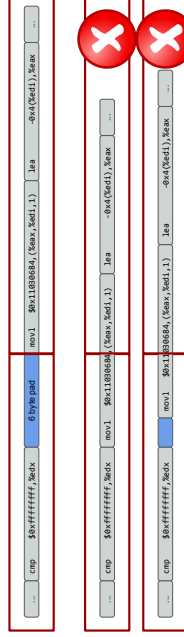
55

Pad Removal



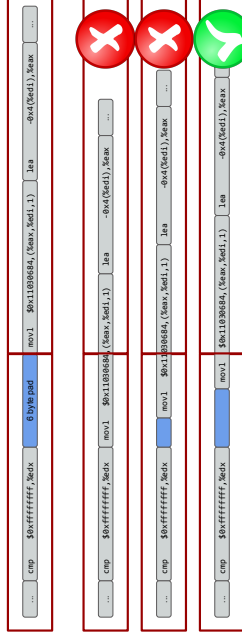
56

Pad Removal



57

Pad Removal



59

NaCI Validator

- One pass: from the start to the end of code
- Maintains two bitmaps: *valid* and *target*
- At each address checks the instruction
- If a valid instruction marks it in *valid* and advance by instruction size
- If indirect branch checks masking instruction presence
- If direct branch, the destination is marked in *target*
- At the end *target* and *valid* are compared together

58

Multipass Validator

- Challenge: Cross-Bundle Instructions

```

1800451: 83 c3 03          or     $0x1,%eax
1800454: 88 41 28          mov    %eax,%eax(%eax)
1800457: 89 51 06          mov    %eax,%eax,%eax
180045c: 8e j2             test  %eax,%eax
1800462: 89 01 00 00      mov    %eax,%eax
1800465: 89 01 00 00      mov    (%%eax),%eax
1800468: 8d 3c 56 00 00 00 lea   0x0(%eax,%eax),%eax1
1800471: 93 49 01          sub   $0x1,%eax4
1800474: 89 49 0f          and   $0xffff,%eax
1800477: 89 49 0f 00 00 00 and   $0xffff,%eax
180047c: 89 49 0f 00 00 00 and   $0xffff,%eax
1800481: c7 6a 38 84 06 03 11 movl  -0x4(%eax),%eax
1800484: 8d 3c 56 00 00 00 lea   0x0(%eax),%eax
1800487: 0f 8a 70 01 00 00      js     180060b
1800490:

```

Multipass Validator

- Challenge: Cross-Bundle Instructions
- Multipass: start validation from every crossing point

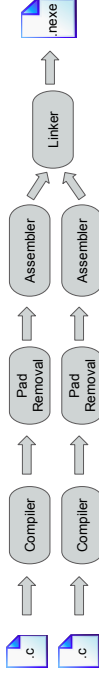
```

1800451: 83 c3 03          or     $0x1,%eax
1800454: 88 41 28          mov    %eax,%eax(%eax)
1800457: 89 51 06          mov    %eax,%eax,%eax
180045c: 8e j2             test  %eax,%eax
1800462: 89 01 00 00      mov    %eax,%eax
1800465: 89 01 00 00      mov    (%%eax),%eax
1800468: 8d 3c 56 00 00 00 lea   0x0(%eax,%eax),%eax1
1800471: 93 49 01          sub   $0x1,%eax4
1800474: 89 49 0f          and   $0xffff,%eax
1800477: 89 49 0f 00 00 00 and   $0xffff,%eax
180047c: 89 49 0f 00 00 00 and   $0xffff,%eax
1800481: c7 6a 38 84 06 03 11 movl  -0x4(%eax),%eax
1800484: 8d 3c 56 00 00 00 lea   0x0(%eax),%eax
1800487: 0f 8a 70 01 00 00      js     180060b
1800490:

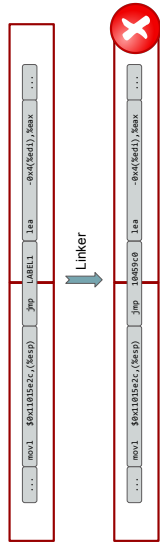
```

Separate Compilation

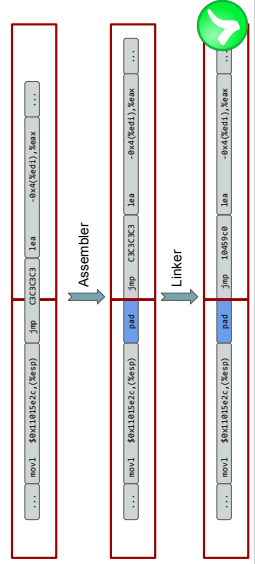
- We process each source file separately
 - Decide about the paddings to be removed
 - Assemble them into object files (using modified GAS)
 - Then link them together



Relocations Problem



Relocations Problem



Thank you
Any Question?