

Vulnerabilities in C/C++ programs – Part II

TDDC90 – Software Security

Ulf Kargén

Department of Computer and Information Science (IDA)

Division for Database and Information Techniques (ADIT)

Integer overflows and sign errors

Adding, subtracting, or multiplying an integer with a too large value can cause it to wrap around

- Can be used to circumvent input validation to e.g. cause buffer overflows

```
void print_user(char* username) {
    char buffer[1024];
    char* prefix = "User: ";
    const unsigned int prefix_len = 6;

    unsigned int len = strlen(username);

    // Space required for prefix, username and
    // string terminator.
    unsigned int size = prefix_len + len + 1;
    if(size > 1024)
        exit_with_error(); // Error, too long string

    strcpy(buffer, prefix); // Copy prefix
    strcat(buffer, username); // Concatenate username

    printf("%s", buffer);
}
```

What happens if the user supplies an extremely long 'username' here?

- If username is longer than `UINT_MAX - 7`, an integer overflow will occur.
- ⇒ Input will pass length check, but still more than 4GB copied into buffer...

Similar problems can arise when casting between data types.

E.g. `int` → `short`:

Most significant two bytes are dropped

Integer overflows and sign errors

A similar class of vulnerabilities are sign errors – mixing signed and unsigned data types in an unsafe way

```
// Reads 'size' bytes from file 'f' into buffer 'out'
void
read_from_file(void* out, FILE* f, unsigned int size);
...

int read_entry(FILE* input)
{
    char buffer[1024];
    int len;

    // Read four-byte length field from file into 'len'
    read_from_file(&len, input, 4);

    if(len > 1024)
        return ERR_CODE; // Error, data won't fit

    // Read 'len' bytes from file into buffer
    read_from_file(buffer, input, len);
    ...
}
```

The problem here is that signed and unsigned data types are mixed.

- What happens if the length field in the file is a negative number, e.g. -1?
 - ⇒ The length check will succeed, as $-1 < 1024$
 - ⇒ In the call to 'read_from_file', the 'len' variable will be interpreted as an unsigned data type
 - ⇒ The 32-bit representation of -1 is $0xFFFFFFFF \approx 4$ billion, way more than the buffer size!

Integer overflows and sign errors

Can be extremely subtle!

If the length check from previous example is changed from this...

```
if(len > 1024)
    return ERR_CODE; // Error, data won't fit
```

... to this, the code is no longer vulnerable. Why?

```
if(len > sizeof(buffer))
    return ERR_CODE; // Error, data won't fit
```

- The value returned by the 'sizeof' operator is always of an unsigned type (size_t)
- According to the C standard, if two values of different data types are compared, and one of the types can represent larger numbers than the other, the value of the smaller type is implicitly cast to the larger.
- The above comparison becomes `if((size_t)len > sizeof(buffer))`
- ... but don't rely on these sort of things to avoid vulnerabilities :-)

Avoiding integer errors

- Again: Perform input validation!
 - Catch e.g. negative lengths of strings, etc.
- Avoid mixing signed and unsigned data types, as well as types of different sizes. **Heed compiler warnings!**
- Understand sizes and conversion rules for data types!
- Use the type 'size_t' for variables representing lengths of things.
 - Is always an unsigned data type (cannot be negative)
 - Guaranteed to be able to represent the length of any object in memory (i.e., it's 32 bits on a 32-bit system and 64 bits on a 64-bit system).
- Check for wraparounds :

```
size_t A = ...
size_t B = ...
if(A > SIZE_MAX - B)
    exit_with_error(); // overflow
size_t sum = A + B;
...
```

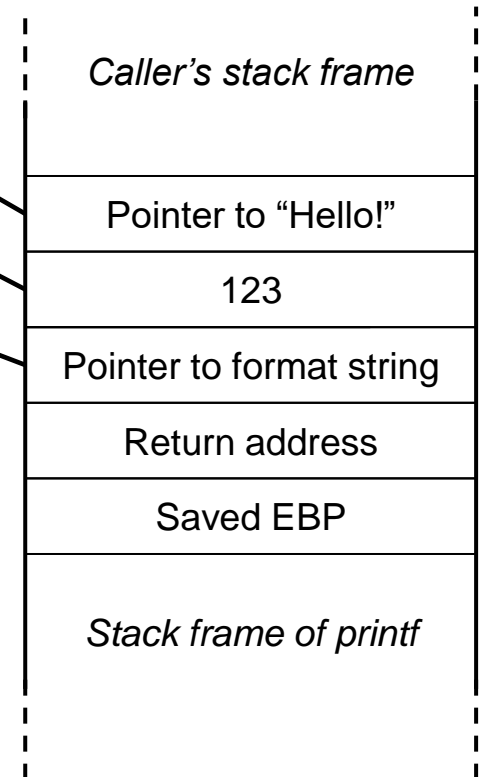
Format string bugs

The printf-family of functions are used in C to format output.

- Takes a format string with placeholders for variable output fields, and a number of arguments corresponding to placeholders in string.

```
printf("An integer: %d, a string: %s", 123, "Hello!");  
// Output: An integer: 123, a string: Hello!
```

- Vulnerability stems from lazy programmers writing `printf(string_from_user)` instead of `printf("%s", string_from_user)`
 - This works fine, as long as the user-controlled string doesn't contain format specifiers!
- `printf` simply assumes that arguments corresponding to all format specifiers exist on the stack – will output whatever is on the stack if that is not the case!
- Supply e.g. a string `"%X%X%X%X"` to output four 32-bit words from callers stack frame in hexadecimal notation – trivial information disclosure.
 - Also possible to read memory at arbitrary address with some trickery.



Format string bugs

- printf also has little known (and used) format specifier %n that is used to store the number of written characters so far into a variable

```
printf("A string: %s%n", "Hello world!", &x);  
// Output: A string: Hello world!  
// x == 22 after execution
```

- Can be used by attacker to write arbitrary data to arbitrary address in memory!
 - E.g. some function pointer at a known address, which is later used for a function call
- Idea (to write arbitrary 32-bit value):
 - Supply the address to write to *in the format string itself*
 - Use a (large) number of format specifiers to advance printf's internal argument pointer to *the format string in the caller's stack frame* (to get to the write address)
 - Control value written by controlling length of string
 - Repeat four times, writing one byte at a time
- Details not important here – available in extra reading material for interested students.

Avoiding format string bugs

- Use `printf("%s", str)` instead of `printf(str)`
 - Unless, perhaps, `str` is a (hardcoded) constant string
- Format string bugs can fairly easily be spotted with static analysis (use of non-constant string as first argument)
- Modern compilers usually warn about (some) insecure use of printf-family of functions.

Summary: Arbitrary Code Execution

Anatomy of an arbitrary code execution exploit:

1. **Supply** executable code (shellcode)
 - a. Inject shellcode into the memory of the process
 - Supply in input strings/buffers
 - Put in environment variable
 - b. Locate shellcode in memory
 - NOP-sled
 - Register trampolines
2. **Redirect** execution to shellcode by **overwriting pointer to code**, which is later dereferenced
 - Return address on stack: **stack-based** buffer overflow
 - Function pointers: **stack/heap-based** buffer overflow, **use-after-free**
 - C++ VTables: **stack/heap-based** buffer overflow, **use-after-free**
 - Also, **format string bugs** allows any pointer with known location to be overwritten

Non-memory-corruption vulnerabilities

So far, we have looked at bugs allowing attackers to overwrite control-data for, e.g., arbitrary code execution or DoS

- Many dangerous types of bugs are not the result of buffer overflows or other memory corruption errors:
 - Race conditions
 - Out-of-bounds reads of data

Race conditions

A shared resource is changed between check and use

```
check_validity_of_user_data()
```

```
[...]
```

```
use_user_data()
```

- Example: File system race conditions

```
if (access(filename, W_OK) == 0) {  
    if ((fd = open(filename, O_WRONLY)) == NULL) {  
        perror(filename);  
        return -1;  
    }  
  
    /* write to the file */  
}
```

- What if file changes between access-check and open?
- Attacker can e.g. replace real file with symbolic link with same name to sensitive file (e.g. /etc/passwd on Unix)

Avoiding race conditions

- Very broad class of vulnerabilities
 - Race conditions on file system
 - Race conditions on memory access between threads
 - etc.
- See literature on course web page for recommendations on avoiding file race conditions in Unix

Out-of-bounds reads

Case study: Heartbleed

Out-of-bounds read from heap-allocated memory in OpenSSL allows attackers to read out certificates, private keys, sensitive documents, etc...

- Due to incorrect implementation of *heartbeat* extension of TLS
- One of the parties in a connection can send a payload with arbitrary data to the other party, which echoes it back unchanged to confirm that it is up and running.
- Problem: Length of payload that is echoed back is not checked. Can read past actual payload into adjacent memory!

Out-of-bounds reads

Case study: Heartbleed

```
int
dtls1_process_heartbeat(SSL *s)
{
    unsigned char *p = &s->s3->rrec.data[0], *pl;
    unsigned short hbtype;
    unsigned int payload;
    unsigned int padding = 16; /* Use minimum padding */
    ...
    /* Read type and payload length first */
    hbtype = *p++;
    n2s(p, payload);
    pl = p;
    ...
}
```

'p' points to data in
SSL record

Record consists of:
Heartbeat type (1 byte)
Payload length (2 bytes)
Payload data (up to 65536 bytes)

'pl' points to
payload data

Copy length of
payload into
'payload'

Out-of-bounds reads

Case study: Heartbleed

```
...
unsigned char *buffer, *bp;
int r;

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

...
/* Enter response type, length and copy payload */
*bp++ = TLS1_HB_RESPONSE;
s2n(payload, bp);
memcpy(bp, p1, payload);
```

Allocate heap memory for reply

Copy 'payload' bytes into buffer for reply message

Problem: The length of 'payload' is never checked! Sender can claim a payload length longer than the actual received SSL record.

- ⇒ Up to 64 kB of adjacent heap memory can be leaked to attacker.
- ⇒ Has been shown to allow reading out private keys from servers!

Avoiding memory safety vulnerabilities

Secure coding practices and principles

- Principles to adhere to
 - Best practices
 - Secure coding standards
- Safer languages

CERT top 10 Secure Coding Practices

1. Validate input
2. Heed compiler warnings
3. Architect and design for security policies
4. Keep it simple
5. Default deny
6. Adhere to the principle of least privilege
7. Sanitize data sent to other systems
8. Practice defense in depth
9. Use effective quality assurance techniques
10. Adopt a secure coding standard

CERT C Secure Coding Standard (excerpt)

Recommendations

- INT01-C: Use `rsize_t` or `size_t` for integer values representing size of an object
- MSC15-C: Do not depend on undefined behavior
- SRC06-C: Do not assume that `strtok()` leaves the parse string unchanged
- FIO07-C: Prefer `fseek()` to `rewind()`
- MEM01-C: Store a new value in pointers immediately after `free()`

Rules

- INT32-C: Ensure that operations on signed integers to not result in overflow
- MSC33-C: Do not pass invalid data to the `asctime()` function
- STR33-C: Size wide character strings correctly
- FIO31-C: Do not open a file that is already open
- MEM32-C: Detect and handle memory allocation errors

... or use a safer language

If not performance-critical:

- Use language with managed memory (Java, C#, Python)
- Note: This generally avoids low-level **memory errors**, but still possible to have serious security problems!
 - Misuse crypto (or other) API:s, insecure deserialization, etc.

Otherwise:

- If “stuck” with C++, consider e.g. using smart pointers (with small performance penalty)
- Consider switching to safer compiled language, like **Rust**
 - Almost as fast as C/C++, however, *security is never “free”* ...

Language design of “safer” languages either

- **impose** (oftentimes unneeded) **safety checks**, with a performance penalty (e.g., Java, C#)
- **prevent you** from designing/implementing code in a *potentially* unsafe (but often convenient) way (e.g., Rust)
 - Harder to learn
 - More upfront work to find suitable abstractions during design

Mitigations

OS and compiler exploit protections

Exploit mitigations

Mitigations are technical measures meant to make attacks *harder*

- Raises cost (time required, expertise) for attackers
- But doesn't necessarily make all attacks impossible

Implemented in either operating system or compiler

- Stack cookies (Compiler based)
- Control-flow integrity (Compiler / Compiler + OS based)
- DEP (OS based)
- ASLR (OS based)

Stack cookies

- Implemented in compiler, must be applied during compilation
- A stack *cookie* or *canary* is inserted in stack frame before the return pointer
- Cookie is checked prior to executing 'ret' instruction. If it has changed, program is terminated with an error message.
 - ⇒ Impossible for attacker to overwrite return pointer with a *buffer overflow* without altering cookie.
- Typical implementation works approximately like this:
 - Cookie placed before saved EBP – prevents overwrite of both return address and saved EBP
 - Cookie stored in global variable that is randomly generated at program startup
 - Static cookies won't work, can just be replicated by attacker!
 - A call to a function that checks cookie integrity is inserted before 'ret' instruction. Terminates program if cookie doesn't match original.
 - Typically also reorders local variables in stack frame so that buffers (arrays) are located first – prevents overwrites of e.g. function pointers in local variables.

Stack cookies

Example

```
void foo(char* input)
{
    // Push global cookie to stack

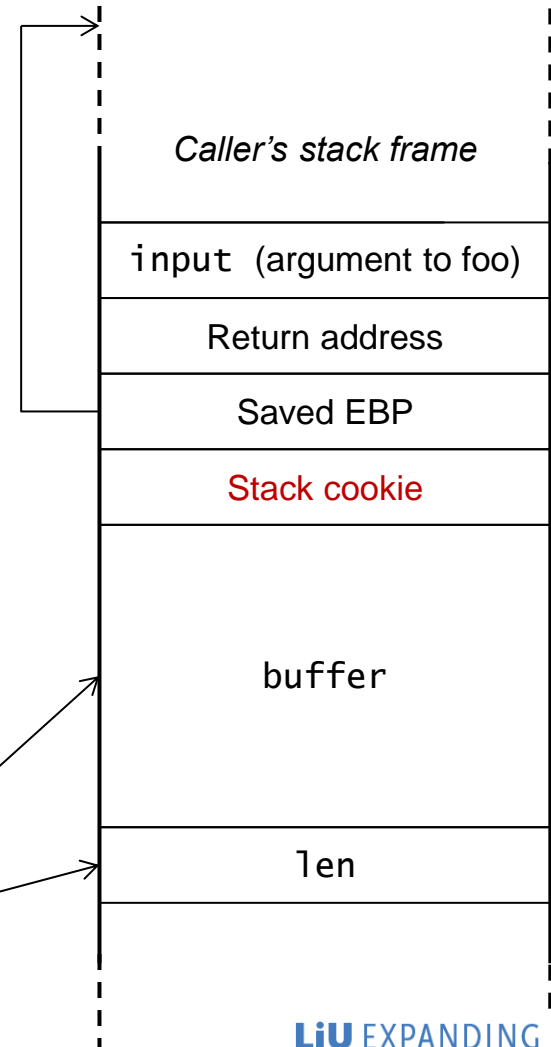
    unsigned int len;
    char buffer[16];

    len = strlen(input);
    strcpy(buffer, input);

    printf("%s: %d\n", buffer, len);

    // Check that cookie match global
    // cookie. Terminate otherwise.
}
```

Note:
Reordered



Defeating stack cookies

- **Only** mitigates stack-based buffer overflows
- Applying stack cookies comes at a cost – for small functions that are called frequently, cost of cookie check can be significant
 - ⇒ Not applied to all functions – various heuristics to determine where to use stack cookies
 - ⇒ Only used in functions with buffers of certain types and sizes – some attacks may still be possible
- On Windows, the Structured Exception Handler (SEH) record on the stack can be overwritten to take control before the return and cookie check

Control-flow integrity (CFI)

- Check at runtime that the target of an indirect branch is valid
- Most commonly used to check that indirect call targets are valid
 - ⇒ Protects against function-pointer overwrites, use-after-free, etc.
- Implemented in e.g. modern Windows versions and (partially) in the LLVM and GCC compilers

Tricky to implement well!

- Need to maintain a *whitelist* of **all** valid targets
 - Backwards compatibility issues (What about legacy libraries that do not have a whitelist?)
 - For example: Still only experimental support for CFI together with legacy libraries in LLVM, not supported at all by GCC
- Checks are made at every indirect call – need to be *very* fast
 - Requires very fast lookups in whitelist...

CFI Example: Microsoft Control Flow Guard

Practical implementation of CFI used in Windows

- Requires support from *both* compiler and OS

Compiler does:

- Store a whitelist of all valid function call targets in generated executable
- Insert calls to a check-function (cf. stack cookies) *before all indirect calls*
 - Takes target function address as parameter
 - By default a check-function that does nothing is used (to make the program runnable on older OS versions)

OS does:

- Creates a *bitmap* of valid addresses for each loaded executable (program or library), using the stored whitelist in executable
 - Each 8-byte unit of memory has an entry in bitmap that says if it contains a valid call target
 - Legacy libraries without CFG have their entire address range marked as valid
- Replaces all calls to the dummy function in loaded executables with a “real” version
 - Does a lookup in bitmap using supplied address – terminates program if not valid

Control Flow Guard Limitations

- CFG uses a *coarse-grained* whitelist to save RAM – granularity is 8 bytes
 - ⇒ Instructions *close* to a valid function start also passes check
- This can be exploited to bypass CFG by using *ROP-gadgets* (soon to be explained) from a function epilogue right *before* a valid function.

Moral of this story: Practical software-based CFI-solutions typically require a tradeoff between thoroughness and memory/computation overhead

- Completely “waterproof” protections are usually too slow to use
 - ⇒ Tradeoffs often enable potential bypass – with some extra effort

Data Execution Prevention

Use hardware-enforced nonexecutable data pages to prevent shellcode from running

Implemented in many different operating systems under different names

- OpenBSD: W^X (Write xor Execute)
- Windows: Data Execution Prevention (DEP)
- Linux: Variants of the PaX MPROTECT patch for Linux kernel

Data Execution Prevention

Recall: Virtual memory divided into pages (typically 4 kB on x86)

- Pages can be marked as Readable, Writable, and Executable
 - ⇒ Write to non-Writable page results in program termination (Segmentation fault)
- Older CPUs (prior to ~2005) didn't have hardware support to enforce the Executable permission
 - ⇒ Possible to execute code from pages marked as non-Executable
- Modern CPUs have this – the NX-bit (for No eXecute)
 - ⇒ Setting all pages for stack, heap, etc. as non-Executable prevents shellcode from executing.
 - ⇒ Effectively mitigates all code execution exploits from previous slides.

Defeating DEP

The return-to-libc attack

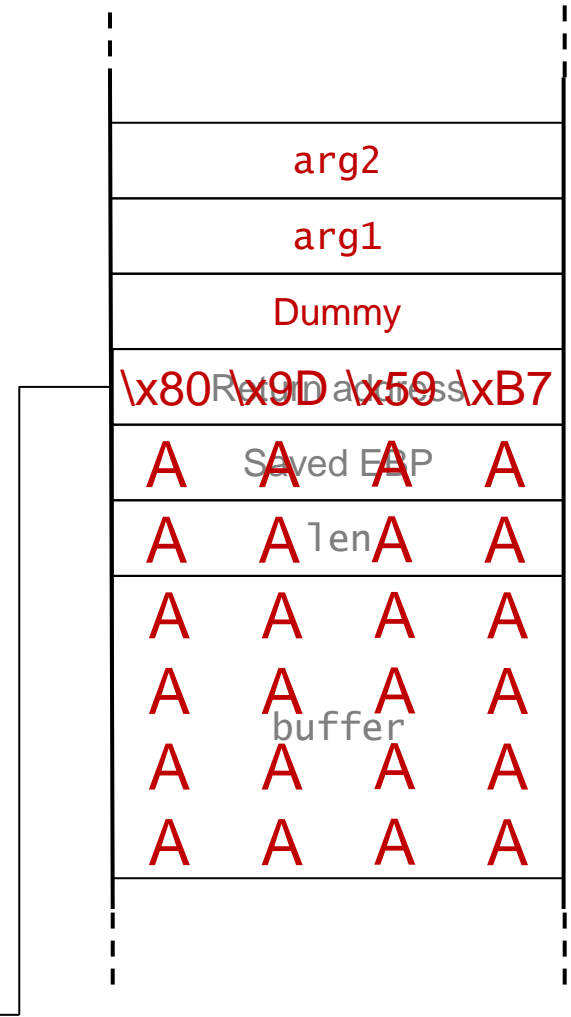
Instead of injecting executable code, *re-use existing function* within program

- Overflow stack buffer to set up stack to look like a function call is about to be made
- Overwrite return pointer to “return” into start of desired function
 - ⇒ No code on the stack is executed – DEP won’t help
- Functions within the standard C library (libc) are popular targets, since libc is present in address space of (almost) every program. Hence the name.
 - E.g. the ‘system’ library function is popular – executes an arbitrary shell command with privileges of calling program

return-to-libc example

Recall the stdcall calling convention:

- Caller pushes arguments from right to left to stack.
- The 'call' instruction pushes return address to stack and jumps to first instruction of called function
- To "call" function `bar(int arg1, int arg2)` using return-to-libc:
 - Overwrite return pointer with address to first instruction of 'bar'
 - Put a dummy value above return pointer. This is where 'bar' expects the caller's 'call' instruction to have put the return address.
 - Put the arguments to 'bar' in correct order on the stack.
 - At 'ret' instruction, 'bar' will be "called", and ESP will point at the dummy "return address", just like in a real call.



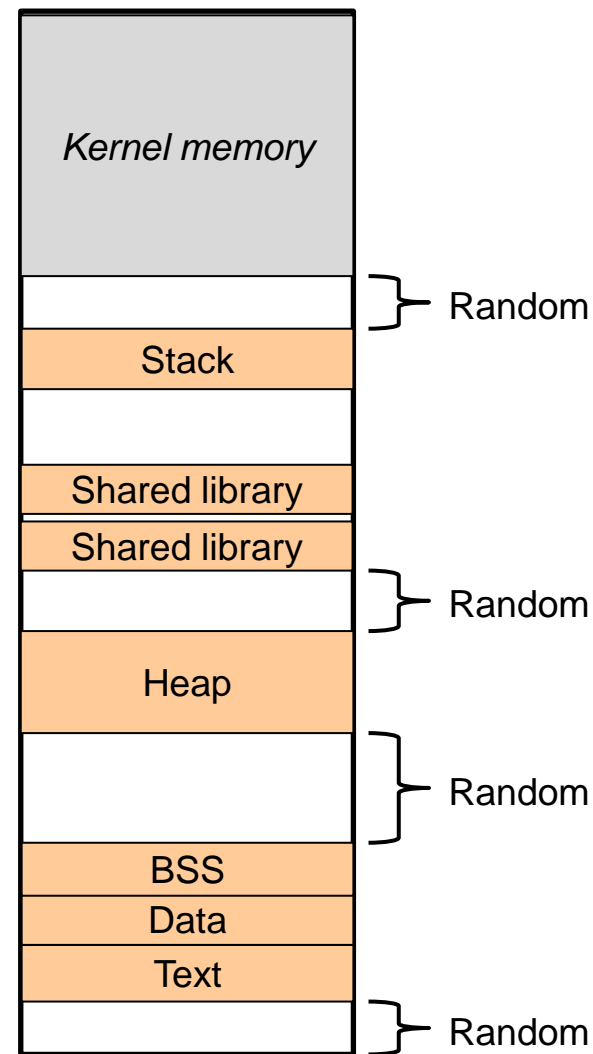
return-to-libc limitations

- Limited to using existing functions within program address space
- Calling functions which takes pointers (e.g. strings) as arguments is tricky.
- Can often not perform calls where one argument is required to have the value zero (Why?)

Address Space Layout Randomization (ASLR)

Observation: Most exploit methods rely on predicting the address of some piece of code or control data.

- Idea: *Randomize* position of heap, stack, main executable, shared libraries, etc. to prevent attacks.
 - New positions each time program is started
- Very effective at mitigating many kinds of attacks.
- Brute forcing still possible on 32-bit machines, where the memory space available for randomization is small. (Works mostly for local exploits.)
- Methods that do not rely on predicting addresses are still effective
 - The *relative* position of data within the same segment is unaffected by ASLR
 - Still possible to e.g. overwrite sensitive non-control data on stack or heap



“Modern” exploit methods

A brief overview

Heap Spraying

Defeats: ASLR

- Applicable in certain scenarios where user controllable input can exert large control over heap allocations
- Make the program allocate large numbers of large memory blocks, filling most of the heap.
 - Each block consists of a large NOP sled followed by shellcode.
- When hijacking control flow of program, e.g. through a stack based-buffer overflow, jump to random position in the middle of the heap
 - ⇒ Large probability of hitting one of the NOP sleds.
- Typically requires a scriptable environment. Popular when e.g. attacking web browsers
 - Create large arrays with e.g. JavaScript, and fill them with NOPs + shellcode.

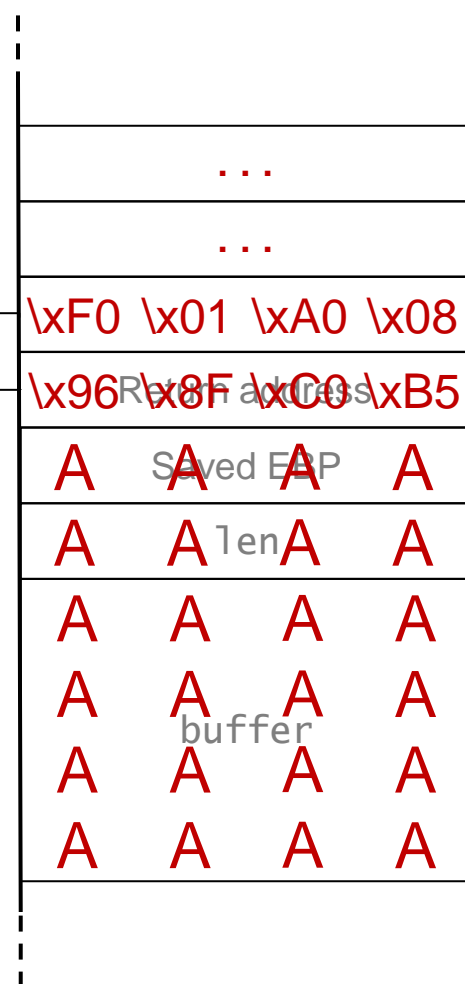
Return Oriented Programming (ROP)

Defeats: DEP

- The “standard” method used today by attackers to bypass DEP
- Generalization of return-to-libc
- First proposed by Hovav Shacham in 2007
 - Showed that a Turing complete “language” could be created by re-using code of an executable.
- Allows *arbitrary* code execution without injecting any code – completely circumvents DEP!
- Idea: Identify code snippets of the form
[do something useful]
ret
in existing code (main executable or libraries).
 - Such snippets are referred to as *gadgets*

Return Oriented Programming (ROP)

- Put addresses of gadgets on the stack, the first one overwriting the return pointer.
 - This “chain” of addresses is often referred to as a ROP chain.
- When the executing function returns, it will pop the gadget address, jump to the gadget, execute the useful instruction(s), and then “return” to the next gadget, and so on.
- Shacham showed that even complex program constructs, such as loops, can be constructed in this way.



```
pop ecx
pop edx
ret
...
```

```
xor eax, eax
ret
...
```

ROP in practice

Most real ROP exploits *pivot* the stack to another attacker-controlled location.

- Allows ROP for non stack-based attacks (function pointer overwrites, use-after-free, etc.)
- Allows for larger ROP-chains

Pivoting principle: Redirect execution to initial pivot-gadget, for example:

```
mov esp, eax  
ret
```

- This will change esp to instead point to whatever eax was pointing to
- eax here points to attacker-controlled part of e.g. heap (Compare to register trampolines!)
- Make sure that ROP chain is set up in memory pointed to by eax. After executing pivot gadget, the "main" ROP chain will start to execute

ROP in practice

Typically not necessary (or possible) to implement your entire shellcode with ROP

- Place regular shellcode at a known address on stack or heap
- Construct a simple ROP chain that just uses, e.g., the `mprotect` system call to *mark the page with the shellcode as executable again*
- Final step of the ROP chain is to simply jump to the regular shellcode

ROP mitigations

ROP attacks rely on being able to predict the addresses of gadgets, and are thus mitigated by **ASLR** – given that the positions of *all* executable memory regions are randomized.

- Still sometimes not the case in practice
 - On older versions of Linux, the executable's section itself ('.text' section) was not randomized, while shared libraries were.
- Some ROP mitigations also target stack pivoting specifically – OS checks that stack pointer actually points to the stack when certain system APIs are invoked. Makes ROP trickier (but not necessarily impossible) in practice.

ROP mitigations

The recently released **Intel CET** technology implements CFI in hardware

- Requires support from both CPU, compiler and OS
 - Needs an 11th generation Intel Core CPU and compatible OS
- CET implements a *shadow stack* – return pointers are replicated on a separate hidden stack by the CPU
- When returning from a function, CPU checks that the return pointer matches the one in the shadow stack
 - Prevents ROP chains from working (in addition to classical return pointer overwrites)
- Also introduces a new ENDBRANCH instruction that is prepended before indirect jump/call *targets*.
 - If the next instruction after an indirect branch is not ENDBRANCH → CPU generates an interrupt and OS kills the process
 - A special whitelist bitmap is needed to allow process to call old libraries without ENDBRANCH instructions
 - Prevents related *call* and *jump* oriented programming attacks

ROP mitigations

Similar implementations of shadow stacks are also offered by other CPU vendors (e.g., AMD, ARM)

Intel has touted CET as the end of ROP exploits

- Remains to be seen if that promise holds true
- Similar things have been said about all of the mitigations mentioned in this course...

Still lots of legacy hardware around...

Effectiveness of mitigations

- No mitigation is a silver bullet
- Some attack methods are thwarted, but often still possible to craft exploits
- However, standard techniques often don't work "out of the box"
 - Often need to combine many different attack techniques, several different vulnerabilities, and program or OS-specific "tricks"
- Example:
 1. Take advantage of a flaw in particular ASLR implementation, or use an information leakage bug, or find target-specific non-randomized executable memory regions to create ROP chain.
 2. Set of gadgets typically limited in practice, create small ROP payload that disables DEP, and jumps to traditional shellcode.
 3. Possibly utilize heap spraying or information leakage bugs to locate shellcode in memory

Effectiveness of mitigations

- Bottom line: Crafting exploits still possible, but requires considerable expertise and time.
 - ⇒ People rarely write exploits “for fun” anymore
- Instead:
 - Professional penetration testers
 - Organized crime
 - Intelligence agencies
- A previously unknown vulnerability (“zero-day”) in popular software with reliable exploit can be worth \$1 000 000 or more...

