

Sicurezza informatica

Federico Zhou

November 16, 2021

Data	Argomento	Materiale didattico	Tipo attivita	Altro
11/03/2021	Low-Level attack: Buffer Overflow (Stack-Based) Part 1	Slide Buffer Overflow Aleph One article on Buffer Overflow Shellcoder Handbook	Didattica frontale	
18/03/2021	Low-Level attack: Buffer Overflow Practical Part 2	GDB Cheat Sheet Peda Repository	Didattica frontale	
25/03/2021	Low-Level attack: Heap Overflow on Metadata	Heap Overflow on Metadata Heap Overflow Repository on Metadata	Didattica frontale	
08/04/2021	Low-Level attack: Use After Free (UAF)	User After Free Vulnerabilities	Didattica frontale	
15/04/2021	Low-level defense: Memory safety and Type safety	Defense against low-level attacks Low-fat pointer technique LowFat Pointer Implementation GDB Input handle	Didattica frontale	
22/04/2021	Low-level defense: Canary, ASLR, DEP, Return Oriented Programming (ROP)	Defense against low-level attacks Return Oriented Programming (ROP)	Didattica frontale	
29/04/2021	Low-level defense: Return Oriented Programming (ROP)	ROP Gadgets Tool Blind ROP paper	Didattica frontale	
06/05/2021	Low-Level defense: Control Flow Integrity (CFI)	Control Flow Integrity (CFI) Control Flow Integrity (CFI) paper	Didattica frontale	
13/05/2021	Program Analysis Testing for Security Purpose	Security Analysis Introduction	Didattica frontale	
13/05/2021	Program Analysis I: Symbolic Execution (Introduction)	Symbolic Execution	Didattica frontale	
20/05/2021	Program Analysis II: Symbolic Execution	Symbolic Execution	Didattica frontale	
20/05/2021	Program Analysis II: Z3 Theorem Prover	Z3 (Microsoft) Theorem Prover Z3 github Documentation Z3 python example Z3 example	Didattica frontale	
27/05/2021	Program Analysis II: Fuzzing	Fuzzing Techniques Klee Tutorial 1 Klee Tutorial video Qsym Hybrid Fuzzer	Didattica frontale	
03/06/2021	Side Channel Attack	Meltdown & Spectre	Didattica frontale	

Buffer overflow

A buffer overflow is a bug that affects low-level code, normally a program would crash but an attacker can alter the situation *to steal informations, corrupt valuable information or run an attacker's code of choice*.

They are relevant because C and C++ are still popular, and buffer overflow attacks still happen regularly. They have a long history and there are many approaches to defend against them.

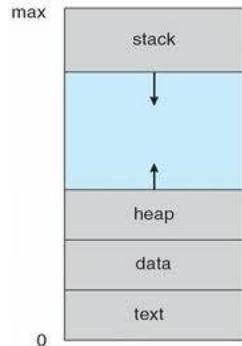
There levels in which we can act to defend against them, the compiler, the OS, or the architecture.

Memory Layout (32bit, Linux)

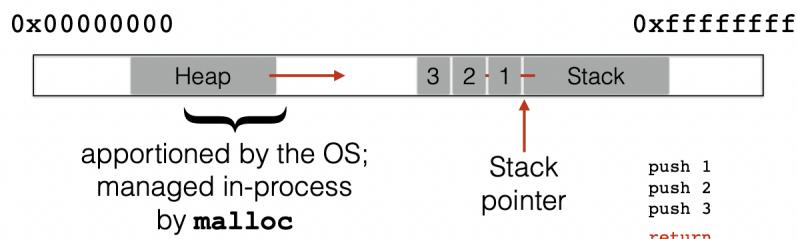
In pile in a 32 bit system, Linux or Intel 32-bit the addresses range from 0x00000000 to 0xFFFFFFFF.

The stack, which grows downwards, is the place in which activation records are stored, together with the heap they are created in runtime. The heap is managed with malloc, calloc,

free and grows upwards, and contains variables. The text segment contains the instructions.



Memory operations make the memory grow upwards, this makes overwriting the SFP possible, allowing memory errors.



Stack and functions: Summary

Calling function:

1. **Push arguments** onto the stack (in reverse)
2. **Push the return address**, i.e., the address of the instruction you want run after control returns to you
3. **Jump to the function's address**

Called function:

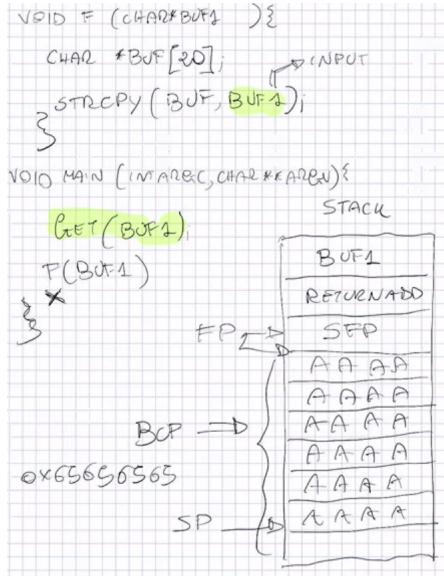
4. **Push the old frame pointer** onto the stack (%ebp)
5. **Set frame pointer** (%ebp) to where the end of the stack is right now (%esp)
6. **Push local variables** onto the stack

Returning function:

7. **Reset the previous stack frame**: %esp = %ebp, %ebp = (%ebp)
8. **Jump back to return address**: %eip = 4(%esp)

Firstly let's talk about buffers, buffers are a contiguous piece of memory associated with a variable or a field, these are common in C. By overflowing a buffer we can put more into the buffer than it can hold.

In this example: BUF1 (which is big [24], is copied into BUF[20], since memory operation increases the heap upwards, it is possible to overwrite the SFP, causing a segmentation fault error.



Now we have crashed the system, but how can we exploit this?

Note that strcpy will let you write as much as you want till a terminator character, we could introduce particularly crafted code to wreak havoc.

If instead of AAAA (which is: 0x65656565) we introduce code which alters the normal execution of the program, overwriting the return address.

How do we build the injection vector?

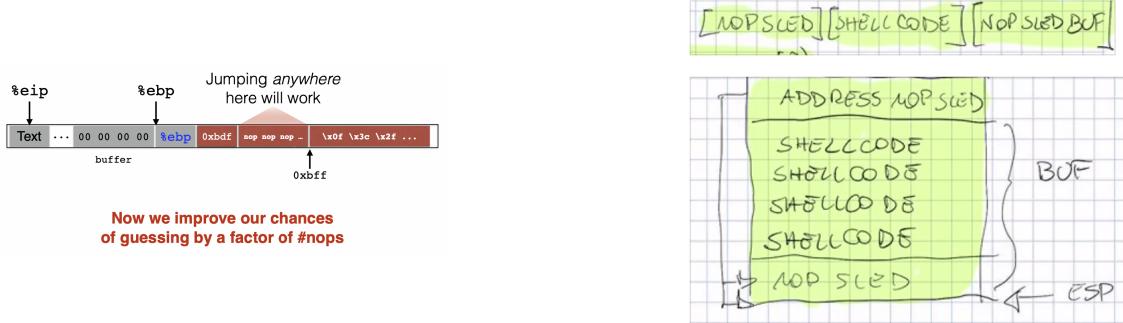


(1) Load my own code into memory

(2) Somehow get %eip to point to it

Note that the "disallocation" only consists in a subtraction of the reference pointer, the data remains. We introduce arbitrary shellcode, and by hijacking the return address we can execute code we want. We can only load machine code into memory, which are instructions already compiled and ready to run, since we're trying to give the attacker general access to the system, the best candidate is the general-purpose shell.

The shellcode is carefully crafted in the format *[SHELLCODE][ADDRESSBUF] = Attack*



Vector and it is copied with the function call "strcpy". It is fundamental we get to know the [ADDRESSBUF], since if we don't an unrecognized instruction is going to run and C is gonna throw an IllegaArgumentException.

Shellcode: machine code -that we want to execute- to run a shell

The idea: we can't insert a "jump to my code" instruction, so we hijack the saved %eip. Since we don't know the desired %eip address we'll have to:

- Try a lot of different values, upwards of 2^{32} possible answers
- Without address randomization, the stack always start from a fixed address, and it doesn't grow very deeply.

How do we increase the chances to execute our shellcodes ?

Introduce the NOPSLED: 0x90 is a NOP instruction, single byte instruction which does nothing. We can fill our vector buffer of NOP operations, thus increasing the chances of execution Depending on the size of the buffer, we have to size our shellcode. The size depends on the allocated space by the target char array in strcpy.

char BUF[20]; strcpy(BUF, BUF1); size of the buffer to work with is 20

The computer by itself has protections, mainly towards the **returnAddress**, such as canary, nx

Heap vs Stack

Let's take an overlook over the differences between the Stack and the Heap:

- The stack:
 - has a fixed memory allocation, which is only known at compile time.
 - inside it reside local variables, return addresses, functions args.
 - it's very fast and automatic, and it is managed by the compiler.
- The Heap:
 - has dynamic memory allocation, known at runtime.

- inside it reside objects, big buffers, structs, persistence and larger objects
- it's slower, manual, managed by the programme

The Heap is a pool of memory used to dynamic allocations, managed with:

- `malloc()`: to allocate memory on the heap
- `free()`: to release memory on the heap

```

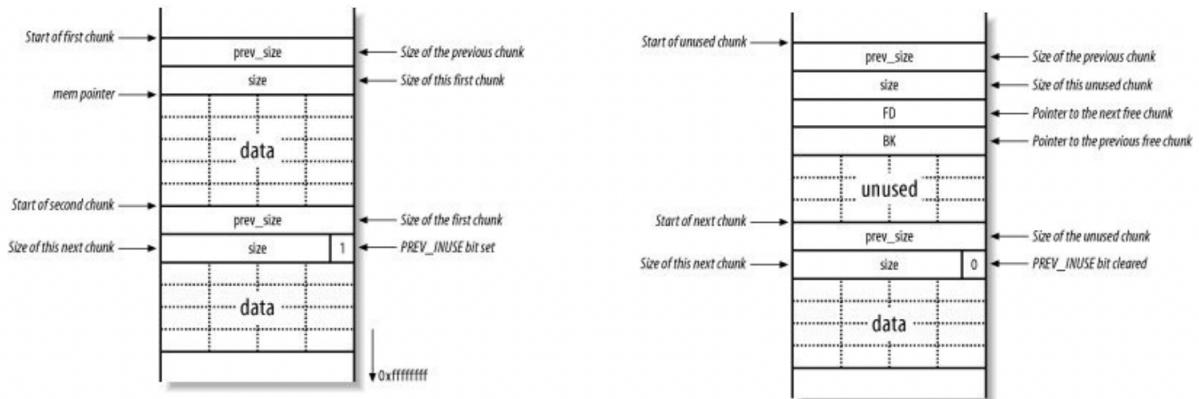
struct malloc_chunk {
    INTERNAL_SIZE_T prev_size; /* Size of previous chunk (if free). */
    INTERNAL_SIZE_T size;      /* Size in bytes, including overhead. */

    struct malloc_chunk* fd;           /* double links -- used only if free. */
    struct malloc_chunk* bk;

    /* Only used for large blocks: pointer to next larger size. */
    struct malloc_chunk* fd_nextsize; /* double links -- used only if free. */
    struct malloc_chunk* bk_nextsize;
};

```

Let's see an overview of a heap chunk:



The chunk is a part of the memory allocated in the heap. Alongside the data, it also contains metadatas:

- `prev_size`: contains the size of the previous chunk
- `size`: it is also equipped with a bit-switch **PREV_INUSE**. It contains the size of the previous chunk.

- FD: the pointer to next free chunk
- BK: the pointer previous free chunk

The heap grows upwards. *prev_size* and *size* concurrently form the metadata section of the heap chunk. When we call $P = \text{malloc}(128)$, the pointer points to memPointer, just below the metadatas.

The free function When $\text{free}(p)$ is used, the LSB of the size field in the metadata of the next chunk is cleared. Additionally the prev_size field of the next size will be set to the size of the chunk we are currently freeing.

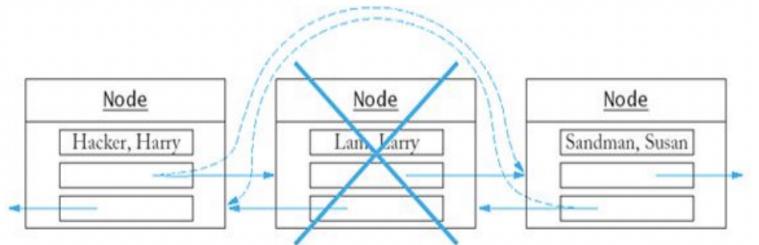
The system keeps multiple lists of the freed chunks, which can be of different size, when a chunk is freed, we can verify if the previous chunk is freed to, and if it is we can unify them into a single bigger chunk.

When an allocation request is made, the system looks into the first free chunk that has a size \geq than the size requested. If no free chunk is found, the top chunk will be used.

AV_TOP è il riferimento al top chunk; rappresenta lo spazio rimanente. Quando solleviamo una memory request e lo spazio nello heap è esaurito eg. $\text{malloc}(256)$ il topchunk è diviso in due, e rispettivamente le parti diventano topchunk e la parte richiesta.

When we deallocate a chunk, we use an unlink function, this leaves an opening for us, the attacker. By hijacking the two metadata values it allows us to write an arbitrary value to an arbitrary location.

```
void unlink(malloc_chunk *P, malloc_chunk *BK,
malloc_chunk *FD)
{
    FD = P->fd;
    BK = P->bk;
    FD->bk = BK;
    BK->fd = FD;
}
```



the intended process

```
...
FD->bk = return address;
FD = P->fd;
BK = P->bk = address of the buffer (injection vector);
FD->bk = BK;
...

```

malicious code which allows us to manipulate the BK

Unfortunately this technique no longer works, the current unlink function verifies beforehand the values We can however look for other flaws:

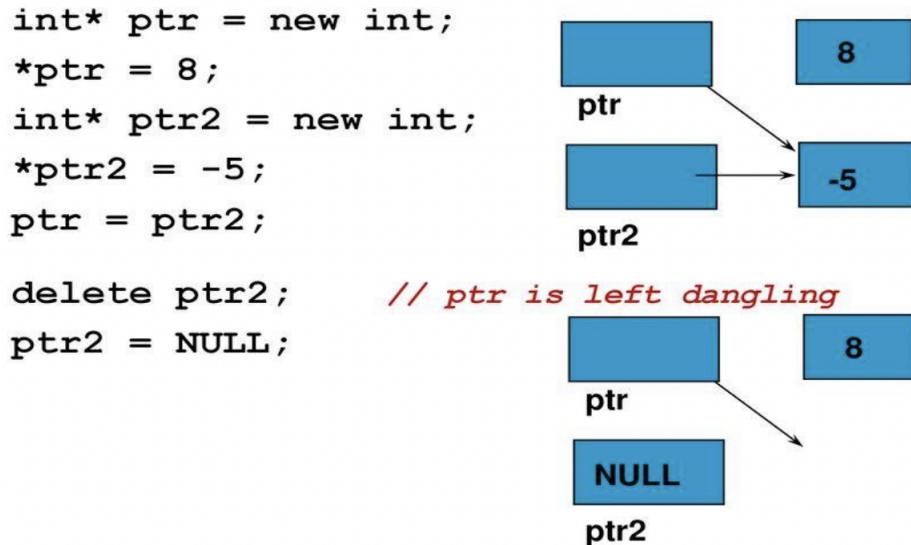
House of Force vulnerability Conditions:

- It requires 3 malloc, more precisely a malloc that allows us to overwrite the top chunk, one malloc with a user controllable size and another call to malloc.
- 1 (or 2) strcpy

Use after free vulnerability Use after free vulnerability are vulnerabilities that happen when a pointer to an object that has been freed is dereferenced. This can lead to information leakage, and furthermore the attacker could also modify unintended memory locations that can potentially lead to code execution.

A prime candidate for an attack is a dangling pointer. Dangling pointer happen when a pointer, which was previously the target of another pointer is deleted, leaving the other pointer dangling.

Leaving a Dangling Pointer



This happens because in C the entire management of pointers is left to programmers. In Java this issue doesn't manifest itself because unused/unallocated memory spaces are collected by the garbage collector.

Let's see an actual example of user after free vulnerability

```

char *retptr() {
    char p ,*q ;
    q = &p ;
    return q ; /* deallocation on the stack */
}

int main( ){
    char *a , *b;
    int i ;
    a = malloc(16) ;
    b = a + 5;
    free(a) ;
    b[2] = 'c' ; /* use after free */

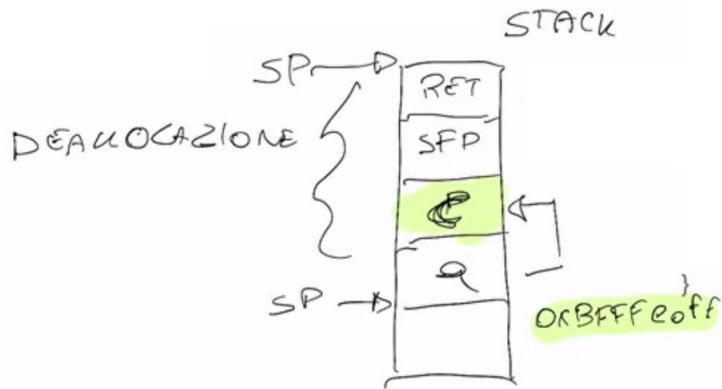
    b = retptr( );
    *b = 'c' ; /* use after free */

}

```

In this example, a is a contiguous memory space long 16 bytes, b points to a location inside of a.

b, after using $b[2] = 'c'$ is able to modify a memory portion of a after it has been freed (de-allocation).



b punta a una porzione dello stack che è stata deallocated. chiamando $*b = 'c'$ modifichiamo una porzione dello stack che è stata deallocated precedentemente (nello stack p diventa c).

Defence and mitigation mechanisms

All these attacks have in common a couple of characteristics:

- The attacker is able to control some data that is used by the program
- The use of data permits unintentional access to some memory area in the program past a buffer, or to an arbitrary position on the stack.

So what can we do to *mitigate* these vulnerabilities?

First of all it's important we mention that we are not going to implement solutions which are going to leave us with a type safety and memory safety language. We are only going to implement mechanisms which are going to make the exploitation of these vulnerabilities too costly/time consuming.

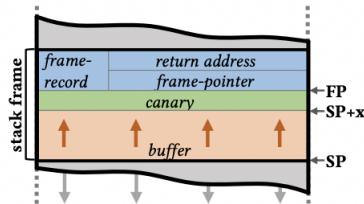
Type safety and memory safety are two fundamental mechanisms which can help us mitigate these vulnerabilities, these properties, if met, ensure an application is immune to memory attacks. Languages such as C and C++ don't implement memory and type safety due to its efficiency load, other languages that do are JAVA or Python.

So what can we actually do?

- Automatic defences:

these defences are specific to the compiler and help protect the return address and - Stack canaries: protects the return address from hijacking.

More precisely a stack canary is a value placed on the stack such that a stack-buffer overflow will overwrite it before corrupting the return address. The buffer overflow can then be detected by verifying the integrity of the canary before performing the return. It is fundamental that the stack canary remains confidential, otherwise the attacker could just copy it and hijack the control flow anyway



- Address space layout randomisation (**ASLR**): is a memory-protection process for operating systems that guards against buffer-overflow attacks by randomising the location where system executables are loaded into memory.

Being able to determine beforehand the position of processes and functions in memory makes exploitation easier, and ASLR is able to put address space targets in unpredictable locations, the attacker won't be able to identify the correct address space, making the application crash and notifying the system.

- Return-oriented programming (**ROP**) attacks are based on the technique that allows the attacker to hijack the program control flow by executing machine code instructions present in the machine's memory.

This is done by manipulating the call stack by exploiting vulnerabilities in a program,

typically a buffer overflow.

Control Flow Integrity (**CFI**) are mechanisms implemented to prevent a variety of attacks based on the redirection of the control flow of the program.

- Secure coding are a set of practises that applies to security considerations that help defend against vulnerabilities, bugs and exploits. Secure coding introduces safeguards which reduce or eliminate the risk of leaving security vulnerabilities in the code, most important of which are **memory safety and type safety**

Memory Safety Memory safety is a fundamental principle which helps mitigate To guarantee a degree of memory safety in C we can:

- only create pointer through standard means

$p = \text{malloc}(\dots)$, or $p = \&x$, or $p = \&buf[5]$ etc.

item only use pointers to access memory that "belong" (that is within the bound) of that pointer and of the type that corresponds. (int* to integers, char* to characters).

This essentially combines two principles: *temporal safety and spatial safety*.

Spatial Safety (in the heap)

Let us Pointers are viewed as triples (p, b, e) where:

- p is the actual pointer
- b is the base of the memory region it may access
- e is the extend (or the bounds) of the region

The defferenced access shall be allowed iff $b \leq p \leq e - \text{sizeof}(\text{typeof}(p))$.

An implementation of a spatial safety compliant memcpy

```
1 void memcpy(void *dst, void *src, int n)
2 {
3     void *dst_base = base(dst);
4     size_t dst_size = size(dst);
5     void *src_base = base(src);
6     size_t src_size = size(src);
7     for (int i = 0; i < n; i++) {
8         void *dst_tmp = dst + i;
9         void *src_tmp = src + i;
10        if (is00B(dst_tmp, dst_base, dst_size))
11            error();
12        if (is00B(src_tmp, src_base, src_size))
13            error();
14        *dst_tmp = *src_tmp;
15    }
16 }
```

Fig. 2. Instrumented version of simple memcpy.

Temporal safety

Memory regions can either be defined or undefined:

- defined means allocated (and active)
 - undefined means allocated, uninitialised, or deallocated
- A temporal safety violation occurs when an attacker tries to access undefined memory space
- spatial safety assures that an accessed region is legal
 - temporal safety assures that region is still in play

When a free function is called, the pointer towards that region is freed up, but the data remains in the stack, just unreferenced. The garbage collection's job is the memory management, its main function is to reclaim memory which was allocated by the program, but is no longer referenced.

The aim is to implement mechanisms similar to the garbage collector, but not quite; we want to protect vulnerable portions of the memory. The garbage collectors are demanding in terms of resources.

But why would we want to use a language that is **not memory safe and type safe**?

The easiest way to avoid all these vulnerabilities would be to use a **memory safe language**, languages that are **type safe are even better**.

Type safety > Memory Safety

The main reason is that C and C++ are **here to stay**, while not memory safe, writing memory safe programs is possible.

The problem of type safety and memory safety has been known for more than 20 years, the limiting factor has always been **performance** (around 12x and 17x overhead). Furthermore adding type safety would also make C and C++ much slower.

Type safety enforcement is expensive, the most important mechanisms are

- **Garbage collection** which avoids temporal violations, but uses much more memory
- **Bound** and null pointer checks which avoids spatial violations
- **Hiding representation** may **inhibit optimisation**
technique such as C-style casts, pointer arithmetic, & operators would not be allowed

So what can we do in C?

- Compiler could add code to **check for violations**
An out-of-bounds access would result in immediate failure just like an *ArrayBoundsException* in Java.

- **No dangling pointers**

Accessing a freed pointer violates temporal safety, accessing uninitialized pointers is also not OK:

```
int *p = malloc(sizeof(int));
*p = 5;
free(p);
printf("%d\n", *p); // violation
```

```
int *p;
*p = 5; // violation
```

- **Type safety**

Each objects shall be ascribed a type (int, pointer to int, pointer to function). Operations on the objects shall be always compatible with the object's type. Type safe programs can't encounter runtime errors. Type safety implies that objects shall always have a compatible types, while guarantee bounds compliance

```
int (*cmp)(char*,char*);
int *p = (int*)malloc(sizeof(int));
*p = 1;
cmp = (int (*)(char*,char*))p;
cmp("hello","bye"); // crash!
```

Memory safe,
but not type safe

So how can we guarantee a degree of type safety in C? Using **enforce invariants**, are properties which cannot be violated by type, by enforcing invariants in the program. Most notably is the enforcement of **abstract types**, which characterised modules. It keeps the **representation of the modules hidden from others**. This guarantees a **degree of isolation** from the rest of the system.

Type-enforced invariants can relate to security properties, by expressing stronger invariants about data's privacy and integrity we can mitigate type mismatch exploits, such as the following in Java.

- **Example: Java with Information Flow (JIF)**

```
int{Alice→Bob} x;
int{Alice→Bob, Chuck} y;
x = y; //OK: policy on x is stronger
y = x; //BAD: policy on y is not
//as strong as x
```

Types have
security labels

Labels define
what information
flows allowed

All of these mechanisms concurrently compose the security by design principle, which means that the language and its capabilities are foundationally secure. In this approach security is built into the system, not mitigating since the design stage notable vulnerabilities.

There are already languages which provide similar features to C, while remaining type safe, notable examples are Google's Go, Mozilla's Rust, Apple's Swift.

Other defensive strategies:

These are complementary, they mitigate bugs but can't shut them out entirely

- Make the bug harder to exploit: introduce more steps, to make them difficult or impossible
- Avoid the bug entirely: practice secure coding practices, advanced code review and testing

To implement these, let's recall the steps of a stack smashing attack:

- Putting the attacker code into the memory. The shellcode must also not contain zero (it is interpreted by the machine as a string terminator)
- Getting `%eip` to point at the shellcode
- Finding the return address (by guessing the raw address)

To make these attacks steps more difficult we can use libraries, compiler mechanisms and the operating system. We are trying to fix the architectural design, and not the code.

Overflows detection using canaries:

Managed by the compiler, canaries work by protecting the return address by detecting any change that could have harmed the stack's integrity. The canaries are placed inside the stack in a protected memory section between user and system variables.

The canary's value has to be random, and unaccessible by the attacker otherwise it could be overwritten with the same value. The canary is introduced when the function is called, using a machine code function called *push canary*

```

    VOID MAIN (argc, argv) {
        CHAR Buf[20]; → PUSH CANARY
        IF (argc < 1)
            PRINTF ("error");
        ELSE
            STRCPY (Buf, argv[1]); → IF [ESP+4] = RANDOM
        }
    }

```

Before calling the *ret* the value of the canary is checked, raising errors if it's not correct. The values the canary can assume are these:

1. **Terminator canaries** (CR, LF, NUL (i.e., 0), -1)
 - Leverages the fact that scanf etc. don't allow these
 2. **Random canaries**
 - Write a new random value @ each process start
 - Save the real value somewhere in memory
 - Must write-protect the stored value
 3. **Random XOR canaries**
 - Same as random canaries
 - But store canary XOR some control info, instead
3. Random XOR canaries work by doing:
1. *[ret XOR canaryValue]* to produce the canary value
 2. *[ret XOR canaryValue]* to verify the canary value

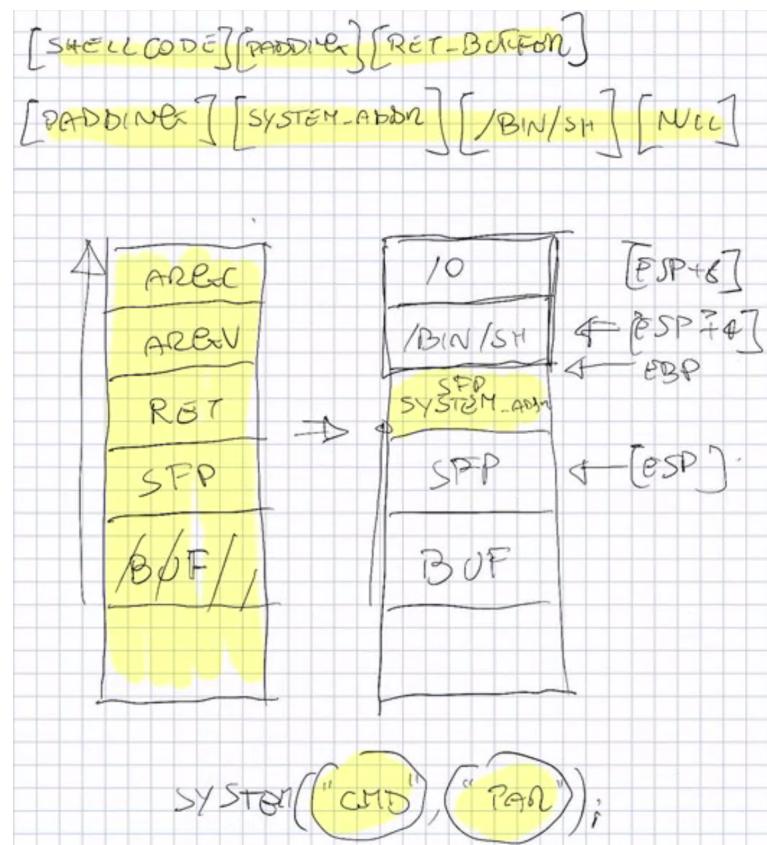
This works due to the nature of XOR, calling it two times gives us the initial value

DEP: Data Execution Prevention:

This defence works by making the stack and heap memory space non-executable. Attackers bypassed this defence by having the shellcode point to system libraries in an attack called *return-to-libc*.

lib-c libraries are used by the gcc compiler, meaning they are always available.

To use this attack, system is called. System is a function that has argument a terminal command. In order to work with this vulnerability, we need to first: call the system function, and two: have it run a meaningful command. This is done by placing the command next to the system call. The injection vector is modified in the following way.



Defence strategies mitigated the *return-to-libc attack* by introducing *Address-space Layout Randomization or ASLR* which places standard libraries and other elements in memory in random places, making them harder to guess. It also has some caveats:

- Only shifts the offset of memory areas, not particular locations within those areas
- May not apply to program code, but only libraries

- Need sufficient randomness, otherwise it could be brute forced.
this makes this technique more promising on 64-bit systems

When we talk about *cat and mouse* we mean the game of cat and mouse being played by attackers and the defence specialists.

Cat and mouse

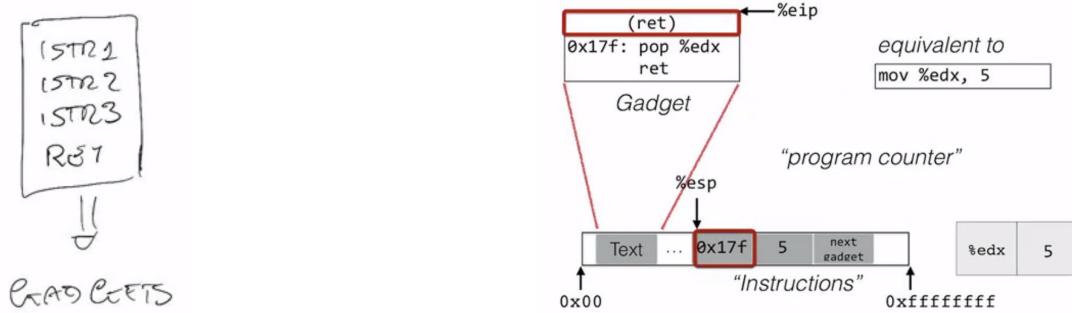
- **Defense:** Make stack/heap nonexecutable to prevent injection of code
 - **Attack response:** Jump/return to libc
- **Defense:** Hide the address of desired libc code or return address using ASLR
 - **Attack response:** Brute force search (for 32-bit systems) or information leak (format string vulnerability)
- **Defense:** Avoid using libc code entirely and use code in the program text instead
 - **Attack response:** Construct needed functionality using return oriented programming (ROP)

Return oriented Programming

First introduced by Hovav Shacham in 2007, the idea is: rather than use a single libc function to run the shellcode, we string together pieces of existing code called gadgets to do it. Intuitively, we need to find the gadgets we need, and we need a mechanism to string them together

gadgets

Gadgets are instruction groups that end with *ret*



The memory structure that we're using is still the stack, which is also where we're loading our code.

- `%esp` serves as program counter
- *gadgets* are invoked via *ret* instruction
- *multiple gadgets* are linked with each others using *ret* instructions
- *gadgets* get their arguments via *pop*, *peek* etc.

In this example we're trying to find an applicable instruction equivalent to *mov %edx, 5*, being: *pop %edx ret*.

So how do I actually find gadgets?

One approach would be an automated search of the target binary for gadgets. This technique work by looking for *ret* instructions, by working backwards. *Are these gadgets actually enough to do anything interesting?*

Yes, Shacham found that with significant codebases (such as libc) gadgets are *Turing complete* meaning they're theoretically capable of solving any computational problem. Schwartz, in 2011, automated the gadget shellcode creation process (ROP compiler), without needing Turing completeness.

In CISC there is a degree of variability in the length of instruction (which is not present in ARM for example since it's a RISC)(RISC instruction lengths are fixed) which provides an extensive degree of freedom to exploits gadgets. This particular aspect is called Dense Instruction Set.

```

push    %ebp          push    %ebx          leave   -1(%eax),%esi add   $0xc,%esp
mov    %esp,%ebp      sub    $0x4,%esp    ret      cmp   $-1,%esi  pop   %ebx
sub    $0x8,%esp      call   8048310     nop      je    8048419  pop   %esi
call   8048304      pop    %ebx          push    %ebp          lea   0x0(%esi),%esi  pop   %edi
call   8048360      add    $0x1c8,%ebx  mov    %esp,%ebp    call   *(%edi,%esi,4)  pop   %ebp
call   8048490      mov    -4(%ebx),%edx  sub    $0x18,%esp  dec   %esi      ret
leave   ret           test   %edx,%edx  and    $-16,%esp  cmp   $-1,%esi  pop   (%esp),%ebx
pushl  pushl 0x80495dc  je    8048326     push    %ebp          lea   0x0(%esi),%esi  not
pushl  pushl 0x80495e0  call   80482d0     add    $0xf,%eax  jne   8048410  mov   %esp,%ebp
add    %al,%eax      pop    %eax          add    $0xf,%eax  lea   0x0(%esi),%esi  not
jmp    *0x80495e4      pop    %ebx          shr   $0x4,%eax  add   $0xc,%esp  push   %ebp
push   $0x0           leave   %ebx        shl   $0x4,%eax  call   80484c4  push   %ebx
push   jmp   8048290     ret      %eax        sub   %eax,%esp  pop    %esi  sub   $0x4,%esp
push   jmp   0x80495e8     nop      %esp        lea   -4(%ebp),%eax  pop    %edi  mov   0x80494f8,%eax
push   push  $0x8           push    %ebp          mov   %eax,0x4(%esp)  pop    %ebp  cmp   $-1,%eax
push   jmp   8048290     mov    %esp,%ebp  movl  $0x80484e8,(%esp)  ret
push   jmp   *0x80495ec  sub    $0x8,%esp  call   80482b0     lea   0x0(%esi),%esi  je
push   push  $0x10          cmpb  $0x0,0x8049600  mov   -4(%ebp),%eax  push   %ebp
push   jmp   8048290     je    804834b     cmp   $0x4d2,%eax  mov   %esp,%ebp  lea   0x0(%esi),%esi
push   jmp   *0x80495f0  add    $0x4,%eax  jne   80483cb    push   %edi  lea   0x0(%edi),%edi
push   push  $0x18          mov    %eax,0x0x0495fc  movl  $0x80484eb,(%esp)  push   %esi  call   *%eax
push   jmp   8048290     call   *%edx,%eax  call   80482a0     mov   %esp,%ebp  mov   -4(%ebx),%eax
push   xor   %ebp,%ebp  mov    0x80495fc,%eax  movl  $0x80484ee,(%esp)  add   $0x119d,%ebx  sub   $0x4,%ebp
push   pop    %esi       mov    (%eax),%edx  call   80482a0     sub   $0xc,%esp  cmp   $-1,%eax
push   mov    %esp,%ecx  test   %edx,%edx  mov   $0x0,%eax  jne   80484b0  pop    %eax
push   and   $-16,%esp  leave   %eax        call   8048278  sub   %ebx
push   push  $0x1           movb  $0x1,0x8049600  ret      -224(%ebx),%eax  pop    %ebp
push   push  %eax         movv  %eax,nop    lea   -224(%ebx),%edx  sub   %ebx
push   push  %esp         leave   %esp        ret
push   push  %edx         ret      %ebp        not
push   push  %ebp         nop      %esp        push   %ebp  push   %ebp
push   push  $0x80483e0    push    %ebp          sar   $0x2,%eax  mov   %eax,-16(%ebp)  push   %ebp
push   push  $0x8048430    mov    %esp,%ebp  push   %edi        je    804847b  mov   %esp,%ebp
push   push  %ecx         sub    $0x8,%esp  push   %esi        xor   %edi,%edi  push   %ebx
push   push  %esi         mov    0x8049508,%eax  push   %ebx        mov   %edx,%esi  sub   $0x4,%esp
push   push  $0x8048384    test   %eax,%eax  pushl 8048483     mov   0x0(%esi),%esi  call   80484d0
call   80482c0      je    8048381    add   $0x11ed,%ebx  lea   0x0(%edi),%esi  pop    %ebx
hlt
nop
nop
push   %ebp         mov    $0x0,%eax  sub   $0xc,%esp  lea   -224(%ebx),%eax  add   $0x1108,%ebx
push   %esp,%ebp      test   %eax,%eax  lea   -224(%ebx),%edx  inc   %esi      call   8048330
push   jmp   8048381    je    8048381    sub   $0x4,%eax  call   *(%esi)  pop    %ecx
push   movl  $0x8049508,(%esp)sar  $0x2,%eax  add   $0x4,%esi  cmp   %edi,-16(%ebp)  cmp   %edi,-16(%ebp)
call   *%eax         jne   8048470  leave   %ebx  lea   0x0(%edi),%esi  not

```

With a simple program like this, we can extract a gigantic amount of instruction.

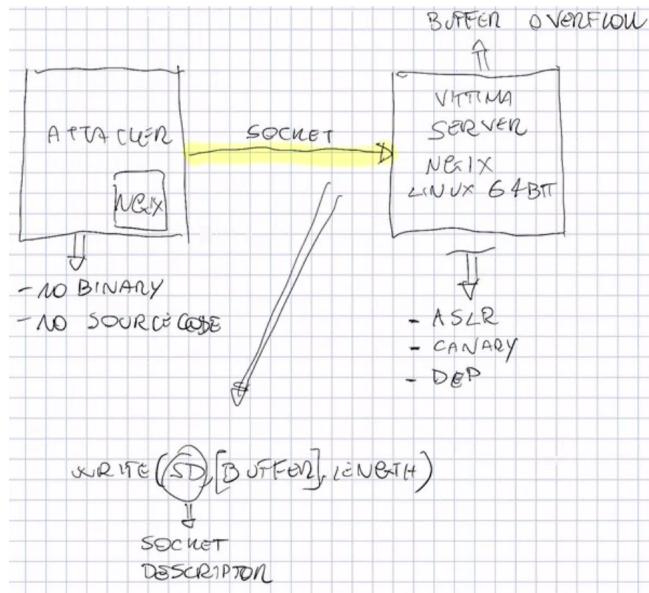
Hacking Blind

Berkley students attack a victim server via a socket, which has all the defence mechanisms we have previously seen enabled, namely ASLR, canary and DEP. The attack happens based on a socket and a linux 64 bit server, via buffer overflow. The exploitation is based on the *write* function:

write(SD, [buffer], length)

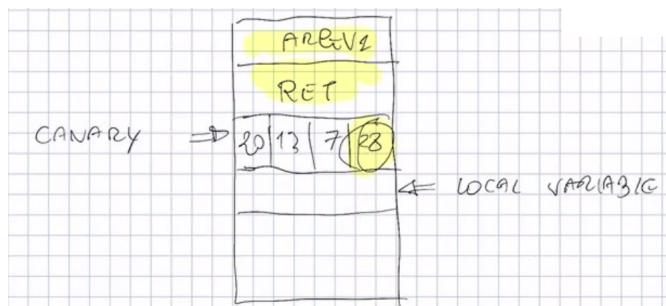
SD is socket descriptor

Utilising a flaw in the *write* function the attacker is able to copy the buffer portion of memory



A fully automated attack based on this exploit yields a shell in under 4000 requests (20 minutes) against a modern system. **Techniques needed:**

1. Generalised stack reading: the possibility to read the stack in its entirety, used to leak the canaries (replace them with the correct values), and to defeat the ASLR (if I am able to read the stack, I am also able to read the location of the program).
Attacking canaries is based on the fact that when a thread crashes on a server, only the thread is reset, not the entire system, and therefore the canary stays the same. The attacker can progressively guess byte per byte the value of the canary. In this example the canary is of size 4 bytes, but on 64 bits system the canary is of size 8 bytes.



To guess the canary we try every value from *0 to 256* every byte, when the thread crashes we proceed with the next value, if the custom value doesn't invoke a crash, we found our value. Using this technique we can also guess the value of the return address.

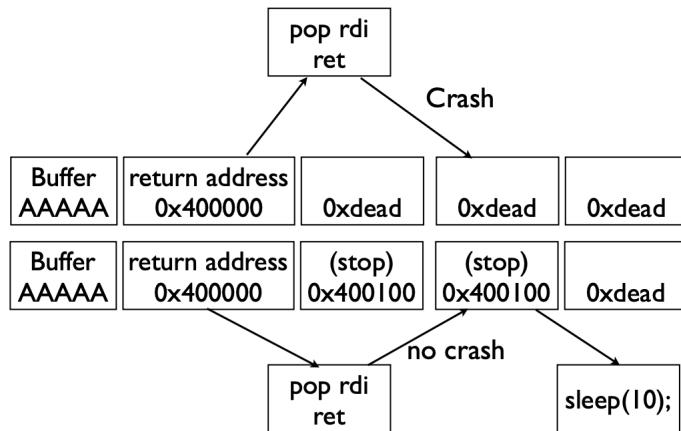
2. Blind ROP: this technique remotely locates ROP gadgets (write functions) The goal is to find enough gadgets to invoke *write*. If we manage this, the binary can be dumped from memory to the network to find more gadgets. As previously mentioned, the *write* system call takes 3 arguments: a socket, a buffer and a length. These are passed in *rdi*, *rsi*, *rdx* registers, and the system call is stored in the *rax* register. The following gadgets are therefore needed:

- 1) pop rdi; ret (socket)
- 2) pop rsi; ret (buffer)
- 3) pop rdx; ret (length)
- 4) pop rax; ret (write syscall number)
- 5) syscall

but how do we find these gadgets?

pop rdi, *rsi*, *rdx* are manageable, *pop rax* and *syscall* are the problematic functions, *syscall* is almost never called since this is usually done through libraries.

To effectively find gadgets, we modify progressively the return address, sorting them into crash gadgets, which cause the socket to crash, non-crash gadgets, which produce an arbitrary effect.



We further elaborate on non-crash gadgets, implementing probes, traps and stops. If the gadget produce the effect we want the stop is invoked and the gadget is selected.

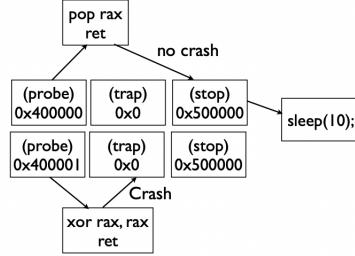


Figure 10. Scanning for pop gadgets. By changing the stack layout, one can fingerprint gadgets that pop words from the stack. For example, if a ‘trap gadget’ is executed rather than popped, the program will crash.

To call functions the PLT (static loader table) and GOT (populated by the dynamic linker) tables are used at runtime to find libraries.

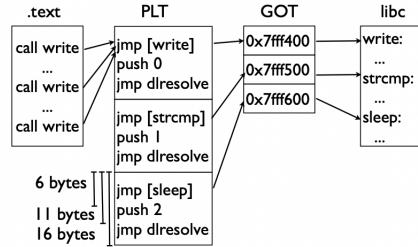


Figure 11. PLT structure and operation. All external calls in a binary go through the PLT. The PLT dereferences the GOT (populated by the dynamic linker) to find the final address to use in a library.

We have seen that memory safety and type safety solve most of the memory errors, *control flow integrity* is the last passive defence mechanism we see.

Control flow integrity

CFI aims to defend a program by observing the program’s behaviour, if the program is not doing what we expect it to do, it might be compromised. To do this we need to define an **expected behaviour**, create mechanisms which let us **detect deviations from expectations effectively**, implement policies which **prevent compromise of the detector**

Control flow graph — CFG

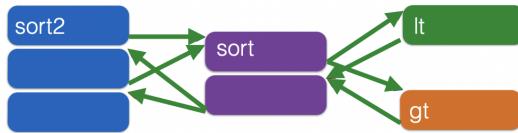
to define expected behaviour. One way we can do this is by utilising **call graphs**. The aim is to represent all the possible combinations of function calls to other functions

```

sort2(int a[], int b[], int len)
{
    sort(a, len, lt);
    sort(b, len, gt);
}

bool lt(int x, int y) {
    return x<y;
}
bool gt(int x, int y) {
    return x>y;
}

```



example of a function which sorts by calling 2 distinct functions

The code is broken up into **basic blocks**, calls are distinguished from **returns**.

In-line reference monitor — IRM

to detect deviations from expectations, it works by introducing control methods inside the program (in binary or inside the program)

Sufficient randomness and immutability

to avoid compromise of the detector

So how efficient is CFI actually?

Classic CFI imposes a 16% overhead on average, 45% in the worst cases, it works on arbitrary executables, and is not modular.

Modular CFI imposes a 5% overhead on average, and 12% in the worst cases, it works only in C (part of the LLVM, source code), and is modular, with separate compilation routines

And how secure is CFI actually?

MFI can eliminate 95% of ROP gadgets on x86-64 versions of SPEC2006 benchmark suites. Average indirect-target reduction is of > 99%, meaning that the percentage of possible targets of indirect jumps is almost completely ruled out.