

Buffer Overflows

Syed Kamran Haider

Department of Electrical & Computer Engineering

University of Connecticut

Email: syed.haider@engr.uconn.edu

Based on and extracted from Nickolai Zeldovitch, Computer System Security, course material at

<http://css.csail.mit.edu/6.858/2014/>



With help from Marten van Dijk

Some of the material is taken from CSE4707: Information Security (Spring'14) by Aggelos Kiayias



Outline

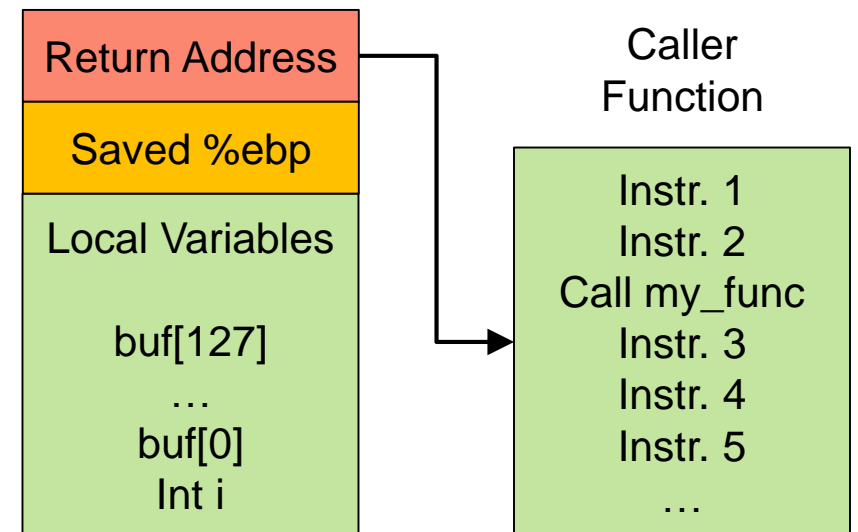
- Introduction of Buffer Overflows
- Payloads of Buffer Overflow Attack
- Avoiding Buffer Overflows
- Mitigating Buffer Overflows
- Advanced Examples
- Detailed Demo

Introduction

- A program execution is broken into several Functions/Procedures
- Before each Function call, its data and some other metadata is placed on the Stack in a **Stack Frame**
- **Return Address** points to the next instruction of the Caller Function to be executed once this function returns.
- **Key Observation: Modifying the Return Address somehow to point to an arbitrary location can be exploited to execute some arbitrary code!**
 - Buffer Overflows

```
void my_func ()  
{  
    char buf[128];  
    int i;  
    gets(buf);  
    // do stuff with buf  
}
```

Stack Frame
of my_func

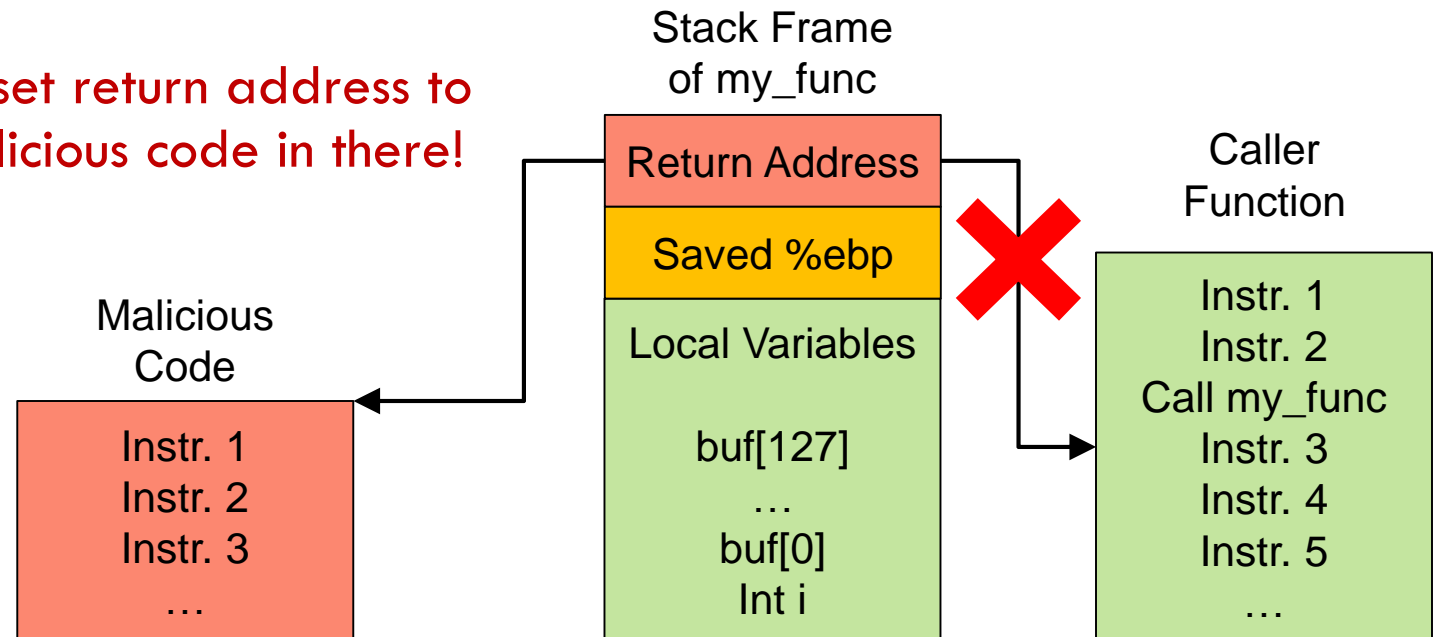


Stack Smashing using Buffer Overflow

How does the adversary take advantage of this code?

- Supply long input, overwrite data on stack past buffer, i.e. create a **Buffer Overflow**!
- Key observation 1: **Attacker can overwrite the return address, make the program jump to a place of the attacker's choosing!**
- Key observation 2: **Attacker can set return address to the buffer itself, include some malicious code in there!**

```
void my_func ()  
{  
    char buf[128];  
    int i;  
    gets(buf);  
    // do stuff with buf  
}
```



Stack Smashing using Buffer Overflow

- How does the adversary know the address of the buffer in the memory?
 - Luckily for the adversary, the Virtual Memory makes things more deterministic!
 - For a given OS and program, Addresses will often be the same...
- Why would programmers write such code?
 - Legacy code wasn't exposed to the internet
 - Programmers were not thinking about security
 - Many standard functions used to be unsafe (strcpy, gets, sprintf)

Payloads of Buffer Overflow Attack

What can the attackers do once they are executing arbitrary code through a Buffer Overflow Attack?

- Use any privileges of the process!
- If the process is running as root or Administrator, it can do whatever it wants on the system.
- Even if the process is not running as root, it can send spam, read files, and interestingly, attack or subvert other machines behind the firewall.

What about the OS?

Why didn't the OS notice that the buffer has been overrun?

- As far as the OS is aware, nothing strange has happened!
- The OS only gets invoked when the application does IO or IPC.
- Other than that, the OS basically sits back and lets the program execute, relying on hardware page tables to prevent processes from tampering with each other's memory.
- However, page table protections don't prevent buffer overruns launched by a process "against itself", since the overflowed buffer and the return address and all of that stuff are inside the process's valid address space.
- OS can, however, make buffer overflows more difficult.

Avoiding Buffer Overflows

1. Avoid bugs in C code
2. Build tools to help programmers find bugs.
3. Use a memory-safe language (JavaScript, C#, Python).

1. Avoid bugs in C code

- Programmer should carefully check sizes of buffers, strings, arrays, etc.
 - Use standard library functions that take buffer sizes into account (`strncpy()` instead of `strcpy()`, `fgets()` instead of `gets()`, etc.).
- Modern versions of gcc and Visual Studio warn you when a program uses unsafe functions like `gets()`.
 - In general, **DO NOT IGNORE COMPILER WARNINGS**. Treat warnings like errors!
- **Good: Avoid problems in the first place!**
- **Bad: It's hard to ensure that code is bug-free, particularly if the code base is large. Also, the application itself may define buffer manipulation functions which do not use `fgets()` or `strcpy()` as primitives.**

2. Build tools to help find bugs

- We can use static analysis to find problems in source code before it is compiled.
- Imagine that you had a function like this:
 - By statically analyzing the control flow, we can tell that “offset” is used without being initialized.
- **Bad: Difficult to prove the complete absence of bugs, esp. for unsafe code like C.**
- **Good: Even partial analysis is useful, since programs should become strictly less buggy.**
 - For example, baggy bounds checking cannot catch all memory errors, but it can detect many important kinds

```
void foo(int *p)
{
    char buf[128];
    int offset;
    int *z = p + offset;
    bar(offset);
}
```

3. Use a memory-safe language

- Use a memory-safe language (JavaScript, C#, Python).
- **Good: Prevents memory corruption errors by**
 - Not exposing raw memory addresses to the programmer, and
 - Automatically handling garbage collection.
- **Bad: Low-level runtime code DOES use raw memory addresses.**
 - So, the runtime code still needs to be correct.
- **Bad: Still have a lot of legacy code in unsafe languages**
 - E.g. FORTRAN and COBOL
- **Bad: Maybe you DO need access to low-level hardware features**
 - E.g., you're writing a device driver.

Why Mitigation?

- All 3 above approaches for “Avoiding” Buffer Overflows are effective and widely used, but buffer overflows are still a problem in practice.
 - Large/complicated legacy code written in C is very prevalent.
 - Even newly written code in C/C++ can have memory errors.
- Therefore, We do need Buffer Overflow Mitigation techniques...

Let's revisit Buffer Overflow Attack!

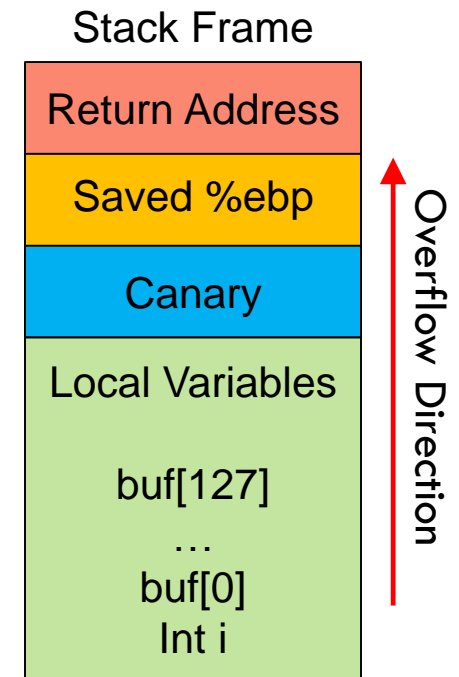
- Two things going on in a "traditional" buffer overflow:
 1. Adversary gains control over execution (program counter).
 2. Adversary executes some malicious code.
- What are the difficulties to these two steps?
 1. Requires overwriting a code pointer (which is later invoked), e.g. Return Address, Function Ptr etc.
 - *Canaries, Bounds Checking etc.*
 2. Requires some interesting/malicious code in process's memory. This is often easier than (1), because it is easy to put code in a buffer because of potentially buggy code!
 - *Non-Executable memory.*
 3. Requires the attacker to put this code in a predictable location, so that he can set the code pointer to point to the evil code!
 - *Address Space Layout Randomization.*

Mitigating Buffer Overflows

1. Canaries (e.g., StackGuard, gcc's SSP)
2. Bounds Checking
 - Electric Fences
 - Fat Pointers (HardBound, SoftBound, iMPX, CHERI)
 - Use shadow data structures to keep track of bounds information (Baggy Bounds).
3. Non-Executable Memory (AMD's NX bit, Windows DEP, W^X, ...)
4. Randomized memory addresses (ASLR, stack randomization, ...)
 - <https://cseweb.ucsd.edu/~hovav/dist/asrandom.pdf>

1. Canaries (e.g., StackGuard, gcc's SSP)

- Idea: OK to overwrite code ptr, as long as we catch it before invocation.
- One of the earlier systems: StackGuard
 - Place a canary on the stack upon entry, check canary value before return.
 - Usually requires source code; compiler inserts canary checks.
- Q: Where is the canary on the stack diagram?
 - A: Canary must go “in front of” return address on the stack, so that any overflow which rewrites return address will also rewrite canary.
- Q: Suppose that the compiler always made the canary 4 bytes of the ‘a’ character. What's wrong with this?
 - A: Adversary can include the appropriate canary value in the buffer overflow!
- Q: Can a Canary protect all buffer overflow attacks?
 - A: No! How about jumping over the canary to overwrite the Return Address?

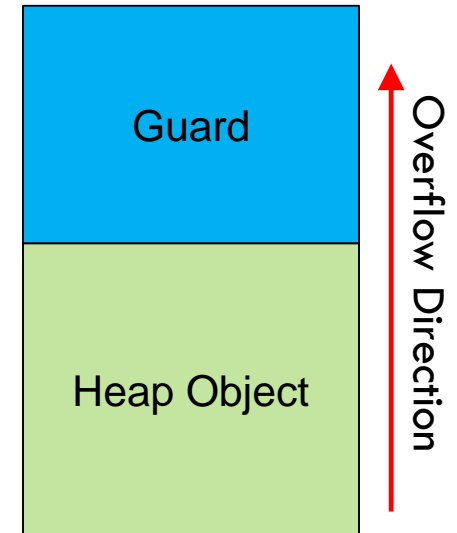


2. Bounds Checking

- Overall goal: Prevent pointer misuse by checking if pointers are in range.
- Challenge: In C, it can be hard to differentiate between a valid pointer and an invalid pointer.
 - E.g. consider an array: `char buf[128];`
 - Consider two pointers: `char *y = buf+100;` `char *z = buf+200;`
 - Which pointer is valid?
- Bounds Checking Goal: For a pointer \underline{y} that is derived from \underline{x} , \underline{y} should only be dereferenced to access the valid memory region that belongs to \underline{x} .

2.1. Electric Fences

- Idea: Align each heap object with a guard page, and use page tables to ensure that accesses to the guard page cause a fault.
- This is a convenient debugging technique, since a heap overflow will immediately cause a crash, as opposed to silently corrupting the heap and causing a failure at some indeterminate time in the future.
- **Big advantage: Works without source code modifications**
- **Big disadvantage: Huge overhead! There's only one object per page, and you have the overhead of a dummy page which isn't used for "real" data.**



2.2. Fat Pointers

- Associate address 'base' and 'bounds' with each pointer.
- Base and bounds checked on each access for security!

Conventional Code

```
void foo(char *str) {  
    int i=0;  
    char buf[16];  
    int x=0;  
    while (str[i] != 0) {  
        buf[i] = str[i]; i++;  
    }  
}
```

Bounds Checked Code

```
void foo(char *str) {  
    int i=0;  
    char buf[16];  
    int x=0;  
    while (str[i] != 0 &&  
           (buf+i > buf.base) &&  
           (buf+i < buf.bound))  
    {  
        buf[i] = str[i]; i++;  
    }  
}
```

- **Problems: Performance Overhead, Incompatibility with existing software!**
- Some recent work: HardBound, Softbound, iMPX, CHERI

3. Non-Executable Memory

- Modern hardware allows specifying read, write, and execute permissions for memory.
- Mark the stack non-executable, so that adversary cannot run their code.
- More generally, some systems enforce “W^X”, meaning all memory is either writable, or executable, but not both. (Of course, it's OK to be neither.)
- **Advantage: Potentially works without any application changes.**
- **Advantage: The hardware is watching you all of the time, unlike the OS.**
- **Disadvantage: Harder to dynamically generate code (esp. with W^X).**
 - Java runtimes, Javascript engines, generate x86 on the fly.
 - Can work around it, by first writing, then changing to executable.

4. Randomized Memory Addresses (ASLR)

- **Observation:** Many attacks use hardcoded addresses in shellcode!
- So, we can make it difficult for the attacker to guess a valid code pointer.
- **Stack randomization:** Move stack to random locations, and/or place padding between stack variables. This makes it more difficult for attackers to determine:
 - Where the return address for the current frame is located
 - Where the attacker's shellcode buffer will be located
- Randomize entire address space (**Address Space Layout Randomization**)
- Can this still be exploited?
 - **Adversary might guess randomness.**
 - On 32-bit machines, there aren't many random bits (e.g., 1 bit belongs to kernel/user mode divide, 12 bits can't be randomized because memory-mapped pages need to be aligned with page boundaries, etc.). [More details: <https://cseweb.ucsd.edu/~hovav/dist/asrandom.pdf>]
- ASLR is more practical on 64-bit machines (easily 32 bits of randomness).

Buffer Overflow Mitigation Summary

Which buffer overflow defenses are used in practice?

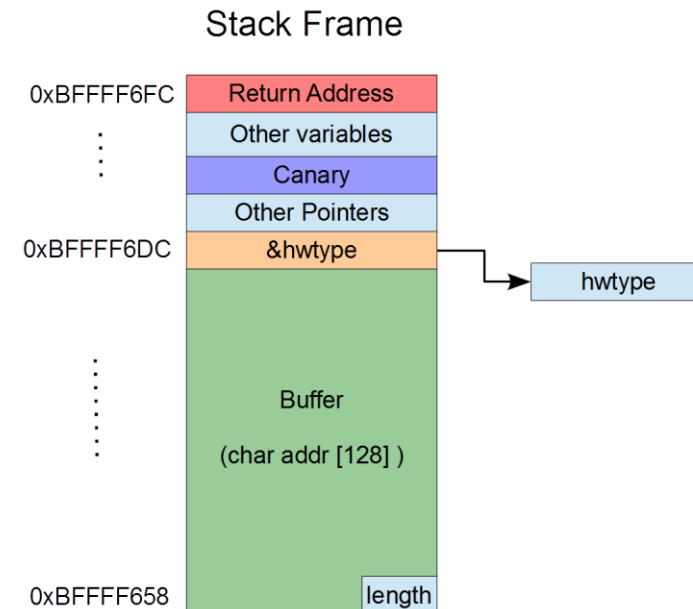
- gcc and Microsoft Visual C enable stack canaries by default.
- Linux and Windows include ASLR and NX by default.
- **Bounds checking is not as common**, due to:
 1. Performance overheads
 2. Need to recompile programs
 3. False alarms: Common theme in security tools: false alarms prevent adoption of tools!
→ Often, zero false alarms with some misses better than zero misses but false alarms.

Ex1: Smashing Stack protected by a Canary

- The program under consideration takes a file as input argument and parses it to print the hardware address stored in the file.
- A structure of type `arp_addr` stores the data read from the file.
- Important members of the structure are
 - `len`
 - `addr[MAX_ADDR_LEN]`
 - `hwtype`
- The correct format of input file is shown below

Type	Length	Address
4 B	4 B	128 B

```
typedef struct{
    ssize_t len;
    char addr[MAX_ADDR_LEN];
    char* hwtype;
    /* Other Members */
} arp_addr;
```



Ex1: Smashing Stack protected by a Canary

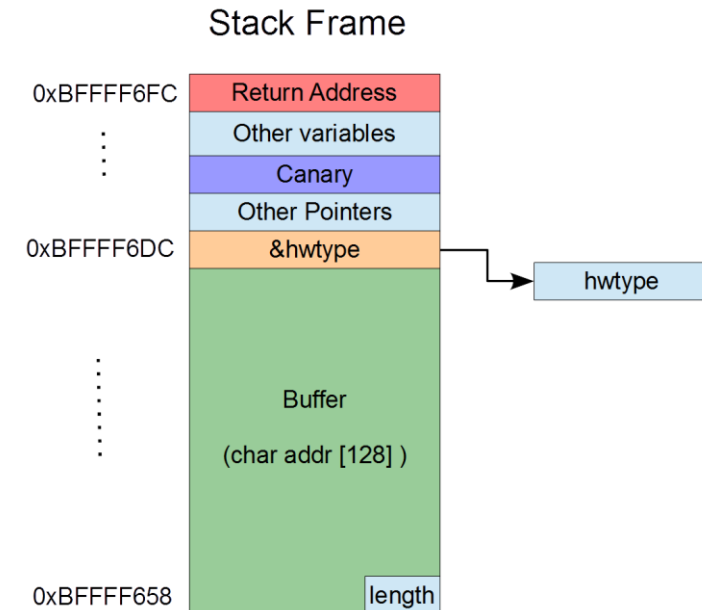
- The `print_address()` function in this program uses a vulnerable function `memcpy()` to copy the input data to the internal data structure.
- First, address length is read from input.
 - Potential to specify incorrect length!
- Specified # of bytes are copied in buffer
 - Possible to overwrite "hwtype"
- 'Type' is stored at location pointed by "hwtype"
 - Possible to overwrite Return address if "hwtype" is pointing to it...!

4 B	4 B	128 B
Type	Length	Address

```

void print_address(char *packet)
{
    arp_addr hwaddr;
    /* Buggy part */
    hwaddr.len = (shsize_t) *(packet + ADDR_LENGTH_OFFSET);
    memcpy(hwaddr.addr, packet + ADDR_OFFSET, hwaddr.len);
    memcpy(hwaddr.hwtype, packet, 4);

    /* Print Address */
    return; }
    
```

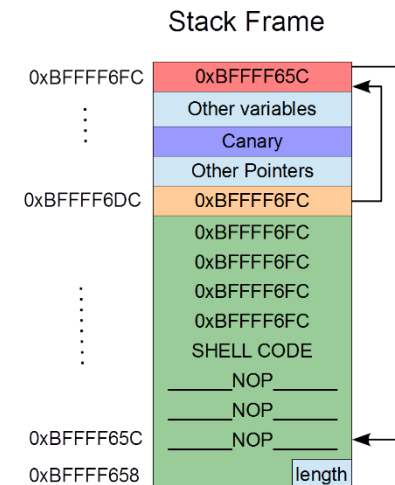
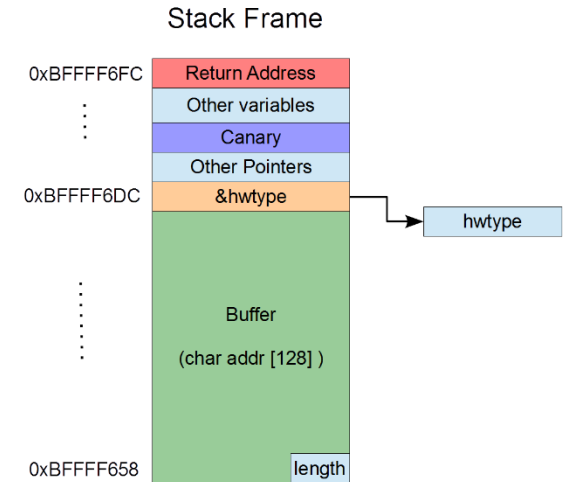


Ex1: Smashing Stack protected by a Canary

- The malicious input stored in file is of the following format

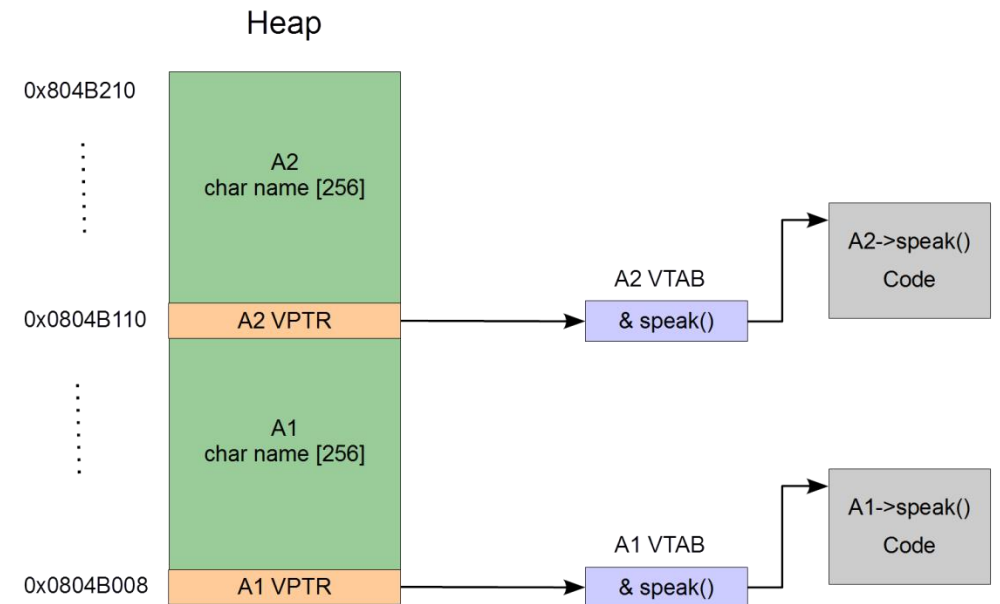


- The buffer overflow **overwrites "hwtype"** but **leaves the canary untouched!**
- "hwtype" now points to the return address
- Writing to the location pointed by "hwtype" basically **overwrites the return address!**



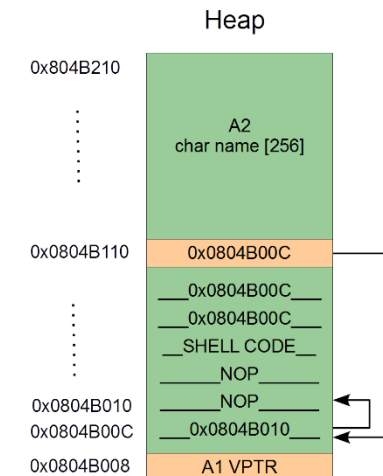
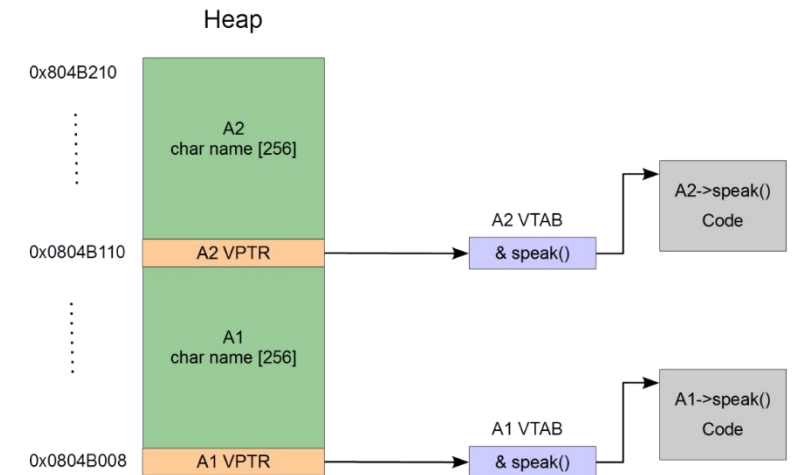
Ex2: Manipulating the Virtual Function Table Pointer (VPTR)

- In Object Oriented Programming, class objects are allocated on heap
 - Stack Smashing Attck cannot work!
- A virtual function of a class can have different definitions for two different objects of the same class.
 - Each class object maintains a *Virtual functions Table* (VTAB) which contains pointers to all the virtual functions of the class and;
 - A pointer to VTAB called *Virtual Pointer* (VPTR) which resides next to the class variables.
 - Depending upon the compiler, VPTR is placed before or after the class variables in the memory.



Ex2: Manipulating the Virtual Function Table Pointer (VPTR)

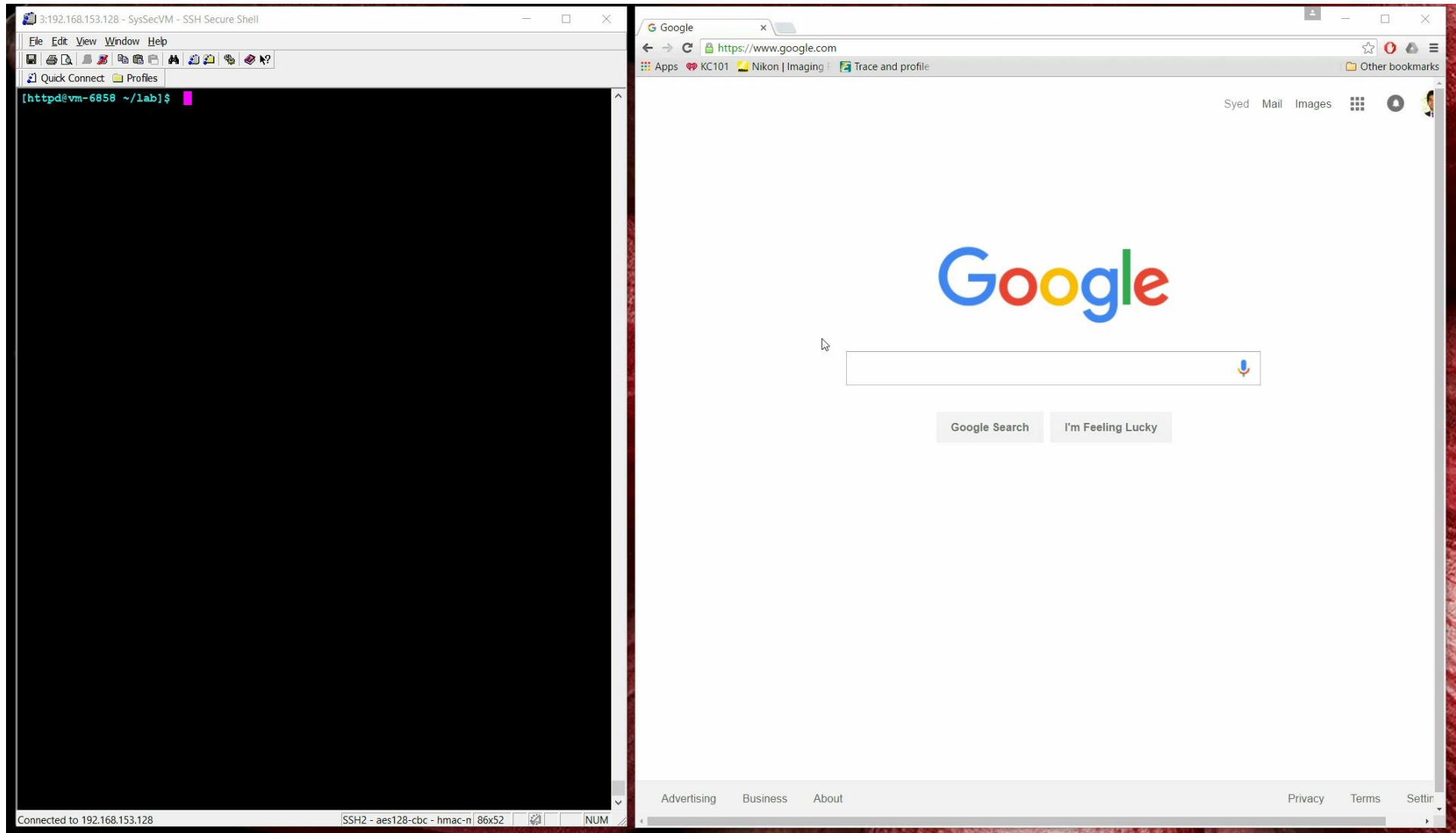
- In this program, name[] buffer of object A1 can be overflown
 - Potential to overwrite VPTR of A2
 - However, we need to be careful while overwriting
- A virtual function call results in two pointer dereferences.
 - First dereference VPTR to go to VTAB
 - Then dereference function pointer stored in VTAB
- We overwrite the buffer as follows:
 - Overwrite A2 VPTR to point to the start of buffer
 - Overwrite the start of buffer with the address of another location in the buffer
 - Handle double dereference
 - Place the malicious code at the second location



Demo: Exploiting Zookbar Website

- Given:
 - **zookws** web server running on a Linux OS
- The zookws web server consists of the following components.
 - **zookld**, a launcher daemon that launches services configured in the file zook.conf.
 - **zookd**, a dispatcher that routes HTTP requests to corresponding services.
 - **zookfs** and other services that may serve static files or execute dynamic scripts.
- Objectives
 1. Identify buffer overflow vulnerabilities in the web server
 2. Crash the web server (Denial of Service requiring a server restart)
 3. Malicious Code Injection (To corrupt/delete some secret files on the server)

Demo: Zoobar Website



Demo: Identifying Buffer Overflows

- The zookws web server consists of the following components.
 - **zookld.c**, a launcher daemon that launches services configured in the file zook.conf.
 - **zookd.c**, a dispatcher that routes HTTP requests to corresponding services.
 - **zookfs.c** and other services that may serve static files or execute dynamic scripts.
- All of these components use functions from the file **http.c**
- **Let's identify a few Buffers that can overflow!**
 1. In HTTP Request Line
 2. HTTP Header Fields

HTTP Request

- An HTTP client sends an HTTP request to a server in the form of a request message which includes following format:
 - A Request-line
 - Zero or more header fields
 - An empty line indicating the end of the header fields
 - Optionally a message-body
- HTTP Request Line example
 - `GET /path_to_file_requested HTTP/1.0`
- HTTP Headers examples
 - `Accept-Language: en-us`
 - `Accept-Encoding: gzip, deflate`
 - `Connection: Keep-Alive`

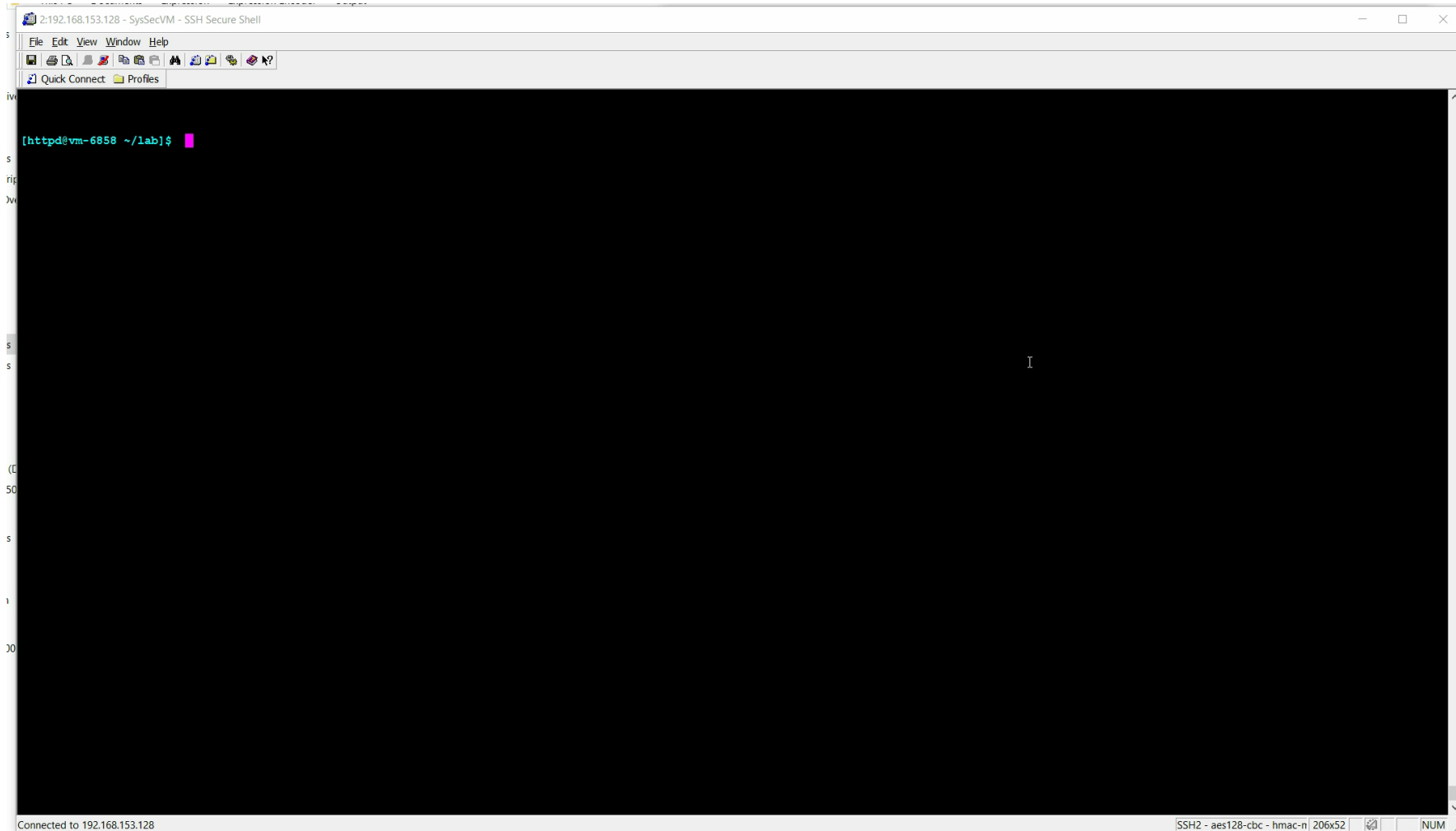
Demo: Identifying Buffer Overflows

Vulnerability in HTTP Request Line's path buffer → zookd.c: reqpath[2048]



Demo: Identifying Buffer Overflows

Vulnerability in HTTP Header's value buffer → zookfs.c: value[512]



Demo: Let's Crash the Web Server

■ Attack Strategy

- Exploit HTTP Request Line Vulnerability, i.e. Buffer `reqpath[2048]`
- Supply an input path that is significantly longer than 2048 bytes → **Overwrite the Return Address**
- Overwriting the return address causes the program to jump to an arbitrary location
- Server crashes due to a segmentation fault!

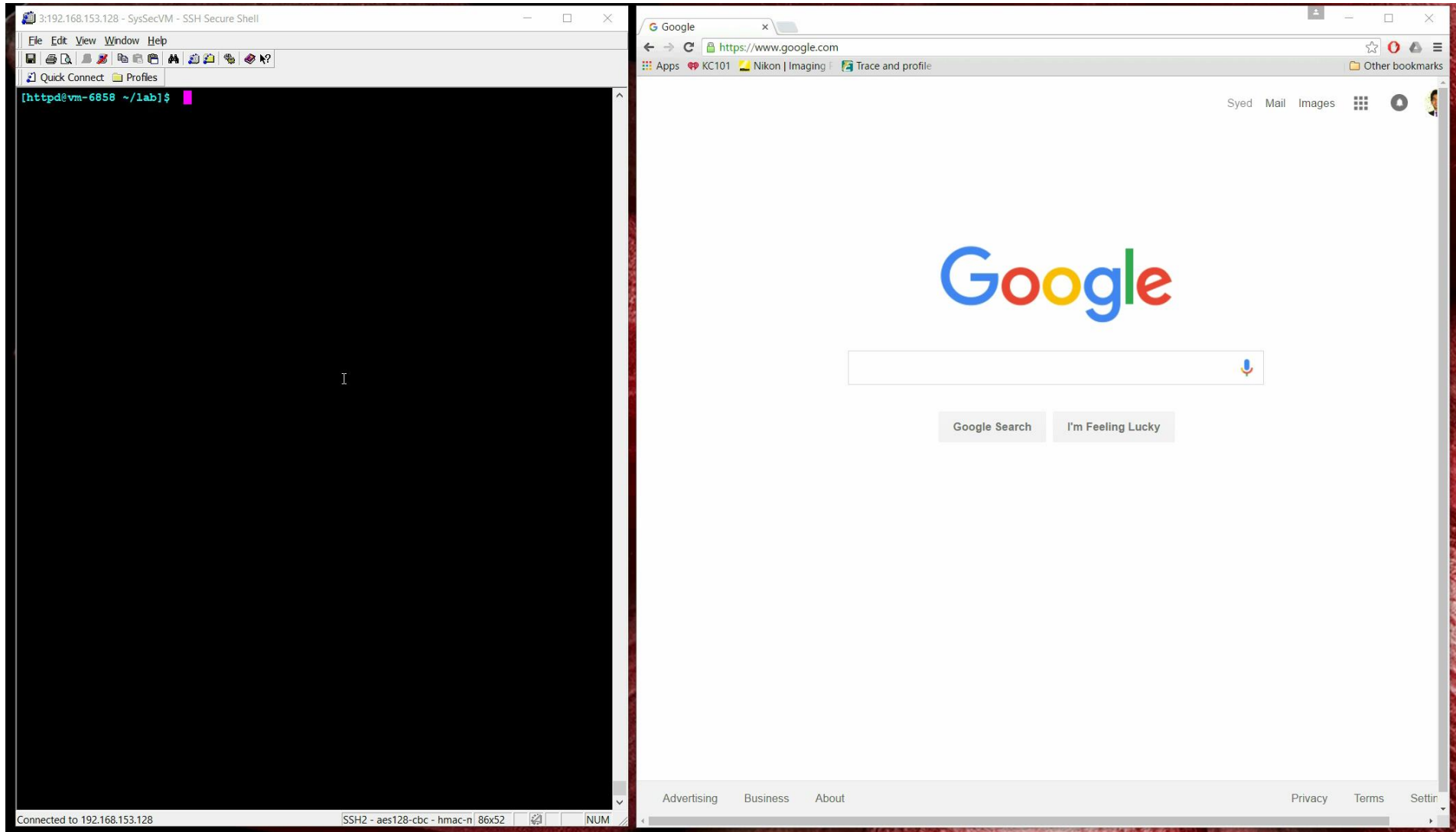
■ Exploit script

- Write a python script to create and issue an HTTP request like the following:
`GET /<very long path> HTTP/1.0`

```
def build_exploit():
    # we want a request line like this: GET /<long path> HTTP/1.0
    req = "GET /"
    for x in range(0, 4096):          # Adding a lot of dummy bytes
        req += "X"
    req += ' '
    req += "HTTP/1.0\r\n" + "\r\n"

    print req
    return req
```

Demo: Let's Crash the Web Server



Demo: Let's Inject Malicious Code

Objective

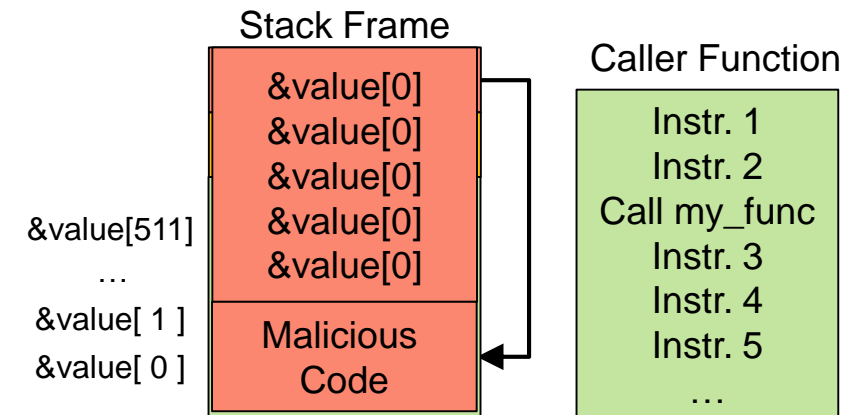
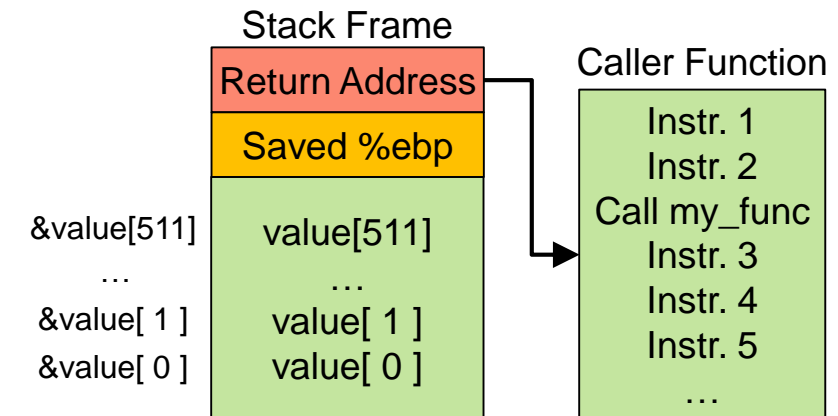
- Delete a file named 'grades.txt' from the server.

Attack Strategy

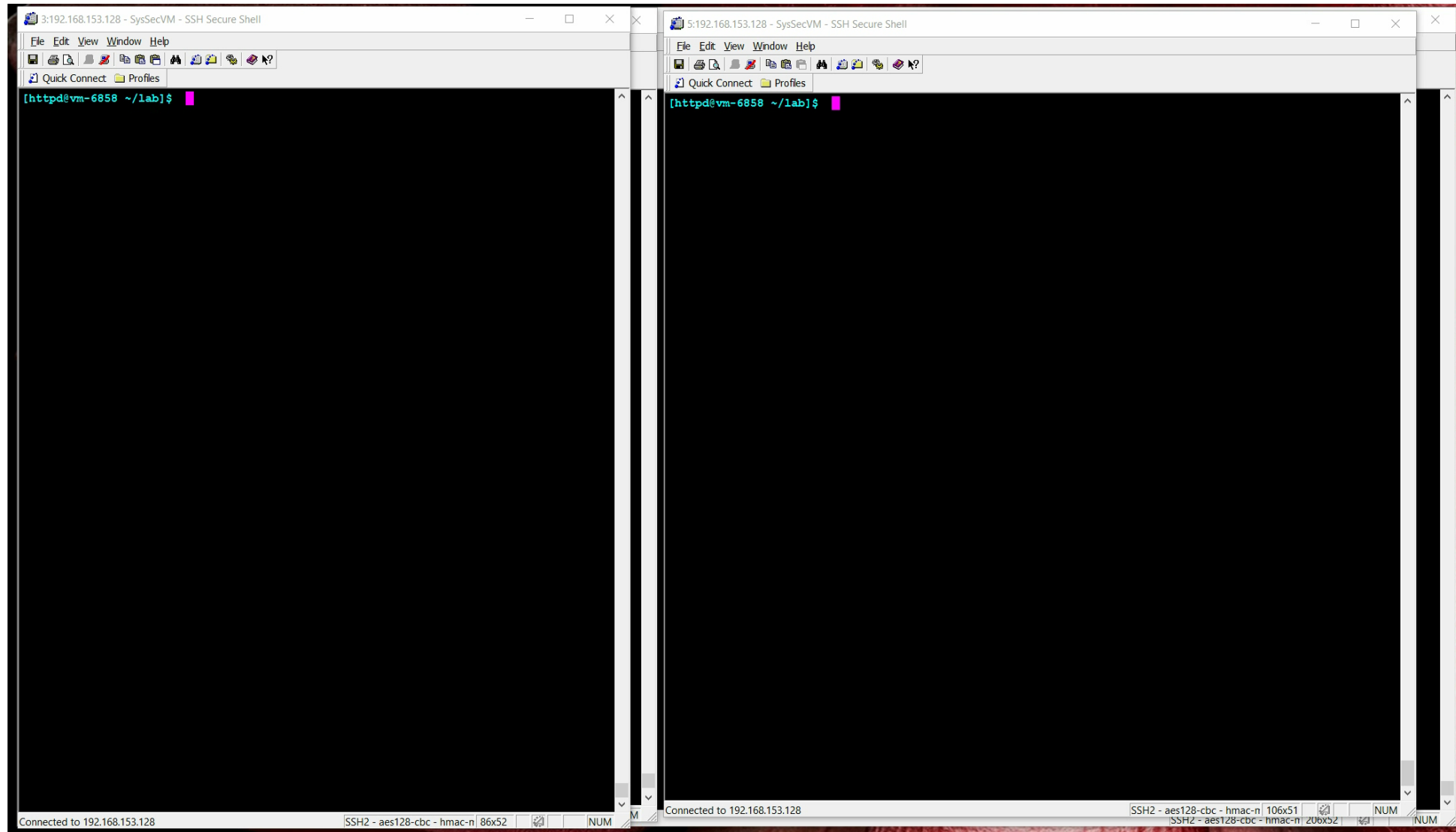
- Exploit HTTP Header Field Vulnerability, i.e. Buffer `value[512]`
- Use GDB to find the address `&value[0]` on stack.
- Prepare a Header value: `<Malicious Code> | <repeat &value[0]>`
- Input this value → **Overwrite the Return Address by `&value[0]`**
- The program jumps to the start of the buffer where malicious code resides.
- Malicious code is executed!

Exploit script

- Write a python script to create and issue an HTTP request like the following:
`GET /<valid_path> HTTP/1.0`
`Accept-Language: <MaliciousCode> | <&value[0]> | ... | <&value[0]>`



Demo: Let's Inject Malicious Code



Thank You!