

Heap Models For Exploit Systems

IEEE Security and Privacy LangSec Workshop 2015

Julien Vanegue

Bloomberg L.P.
New York, USA.

May 24, 2015



Big picture : The Automated Exploitation Grand Challenge

- ▶ A **Security Exploit** is a program taking advantage of another program's vulnerability to allow untrusted code execution or obtention of secret information.
- ▶ **Automated Exploitation** is the ability for a computer to generate an exploit without human interaction.
- ▶ **The Automated Exploitation Grand Challenge** is a list of core problems in Automated Exploitation. Most (all?) problems are unsolved today for real-world cases.
- ▶ Problems relate to: Exploit Specification, Input Generation, State Space Representation, Concurrency Exploration, Privilege Inference, etc.
- ▶ The complete challenge is described at:
http://openwall.info/wiki/_media/people/jvanegue/files/aegc_vanegue.pdf

Today's topic: Heap layout prediction - AEGC Problem I

Disclaimer: this is work in progress research.

Tooling is still in development (no evaluation provided).

Presentation acts on a simplified heap.

Heap can be **non-deterministic**, we focus here on the **deterministic** heap behavior only.

Why is this an important problem?

- ▶ Nowadays, heap-based security exploits are common intrusion software.
- ▶ Exploit mitigations have made writing these exploits an expert's job.
- ▶ Heap allocator implementations are vastly different across Operating Systems.
- ▶ **There is close to no formal research on the topic.**
- ▶ Agenda: Craft and formalize a generic heap exploit technique.

Reminder: Heap vulnerability classes

- ▶ **Heap-based buffer overflow** - Overwrite adjacent memory chunk.
- ▶ **Double free / Invalid free** - Free data that is not a valid allocated chunk.
- ▶ **Use-after-free** - A pointer that was freed is cached and incorrectly used.
- ▶ **Information disclosures** - An attacker can read the content of memory.

Reminder: Heap-based buffer overflow

```
1: char* do_strdup(char *input, unsigned short len) {
2:   unsigned short size = len + 1; // May overflow short capacity
3:   char *ptr = malloc(size); // allocate small amount of memory
4:   if (ptr == NULL)
5:     return (NULL);
6:   memcpy(ptr, input, len); // Buffer overflow may happen
7:   return ptr;
8: }
```

Reminder: Invalid free

```
1: int posnum2str(int x) {
2: char *result;
3: if (x ≤ 0) goto end; // Early exit
4: result = calloc(20, 1);
5: if (result == NULL)
6:     return (NULL);
7: if (num2str(result, x) == 0)
8:     return (result);
9: end: free(result); // May free uninitialized pointer
10: return (NULL);
11: }
```

Reminder: Use-after-free

```
1: char *compute(int sz) {  
2: char *ptr = malloc(sz);  
3: if (ptr == NULL) return (NULL);  
4: int len = f(ptr); // Assume f will free ptr under some conditions  
5: ptr[len] = 0x00; // ptr was already freed!  
6: return (ptr);  
7: }
```

Reminder: Information disclosure

Require: `sock` : Valid network socket

Ensure: **True** on success, **False** on failure

```
1: char buff[MAX_SIZE]
2: int readlen = recv(sock, buff, MAX_SIZE);
3: if (readlen ≤ 0) return False;
4: rec_t *hdr = (rec_t *) buff;
5: char *out = malloc(sizeof(rec_t) + hdr->len);
6: if (NULL == out) return (false);
7: memcpy(out, buff + sizeof(rec_t), hdr->len); // Read out of bound
8: out->len = hdr->len;
9: send(sock, out, hdr->len + sizeof(rec_t)); // Send memory to attacker
10: free(out);
11: return True
```

Original AEGC problem I harness test

```
1: struct s1 { int *ptr; } *p1a = NULL, *p1b = NULL, *p1c = NULL;
2: struct s2 { int authenticated; } *p2 = NULL;
3: F() {
4: p1a = (struct s1*) calloc(sizeof(struct s1), 1);
5: p1b = (struct s1*) calloc(sizeof(struct s1), 1);
6: p1c = (struct s1*) calloc(sizeof(struct s1), 1);
7: }
8: G() { p2 = (struct s2*) calloc(sizeof(struct s2), 1); }
9: H() { free(p1b); }
10: I() { memset(p1a, 0x01, 32); } // Buffer overflow
11: J() { if (p2 && p2->authenticated) puts("you win"); } // Go here
12: K() { if (p1a && p1a->ptr) *(p1a->ptr) = 0x42; } // Avoid crash
```

Goal: Automate heap walk = { F(); H(); G(); I(); J(); }

What do these vulnerabilities have in common?

- ▶ In heap overflow case, attacker expects to place an interesting chunk after the overflowed chunk.
- ▶ In use-after-free case, attacker expects to place controlled chunk in freed memory before it is used incorrectly.
- ▶ In invalid free case, attacker expects to place controlled heap memory at location of invalid free.
- ▶ In information disclosure, attacker expects to place secret in heap just after chunk allowing disclosure.
- ▶ In harness test of Problem 1 (previous slide), we expect chunk p2 to be reusing p1b's memory after it was freed.
- ▶ Summing up: Exploitation depends on location of chunks relative to each others.
- ▶ What is a good **layout abstraction** for the heap?

Studied allocators

- ▶ Doug Lea's malloc (DLMalloc) - Linux.
- ▶ PTMalloc (DLMalloc + thread support) - Linux.
- ▶ Windows heap (including Low Fragmentation Heap).
- ▶ NOT studied: JEmalloc (FreeBSD / NetBSD / Firefox).
- ▶ NOT studied: Garbage Collection (Sweep and Mark algorithm etc).

Typical (simplified) heap allocation algorithm

1. Try to use one of the cached (last freed) chunks.
2. Try to find a fitting chunk in the free chunks list.
3. Try to coalesce two free chunks from free list.
4. If still fails, try (2,3) with each free list in increasing order.
5. If everything fails, try to extend the heap.
6. Otherwise, return an error (NULL)

Formal heap definition

$\mathcal{H} = (\mathcal{L}, \Gamma_a, \Gamma_f, ADJ, Top)$ where:

- ▶ $\mathcal{L} = (l_1, l_2, \dots, l_n)$ is a set of lists of available memory chunks. Each list holds free chunks for a given size range.
- ▶ $l = (c_1, c_2, \dots, c_n)$ are individual memory chunks in list l .
- ▶ $\Gamma_a : l \rightarrow int$ is a **counter map of allocated chunks** for a given size range.
- ▶ $\Gamma_f : l \rightarrow int$ is a **counter map of free chunks** for a given size range.
- ▶ $ADJ : c \times c \rightarrow \mathcal{B}$ is the **adjacency** predicate (true if chunks are immediately adjacent).
- ▶ Top is the current chunk in \mathcal{H} with the highest address.

Heap semantics

Heap primitives:

(F)ree : A memory chunk is freed.

(R)ealloc : A memory chunk is extended.

(A)lloc : A memory chunk is allocated.

(C)oalesce : Two memory chunks are merged.

(S)plit : A big memory chunk is split into two smaller ones.

(E)xtend : The heap is extended by a desired size

Heap transition system:

$$\begin{aligned} \mathcal{H}' &\longleftarrow F \ p \ \mathcal{H} \\ (\mathcal{H}', p_2) &\longleftarrow R \ p_1 \ sz \ \mathcal{H} \\ (\mathcal{H}', p) &\longleftarrow A \ sz \ \mathcal{H} \\ (\mathcal{H}', p_3) &\longleftarrow C \ p_1 \ p_2 \ \mathcal{H} \\ (\mathcal{H}', p_2, p_3) &\longleftarrow S \ p_1 \ off \ \mathcal{H} \\ (\mathcal{H}', p) &\longleftarrow E \ sz \ \mathcal{H} \end{aligned}$$

Key ideas

1. There are two levels of semantics: physical and logical:
 - ▶ The **physical semantic** is concerned with the adjacency of chunks in memory.
 - ▶ The **logical semantic** is concerned with the population of chunk lists.
 - ▶ **Our goal is to reconcile physical and logical heap semantics.**
2. Heap primitives must include user interactions (F, R, A).
3. Core internal heap mechanisms are defined as first class primitives (C, S, E).
4. An Adjacency predicate **ADJ** (used in S and E only) defines the physical semantic. Everything else is house cleaning and defines the logical semantic using two counters per list.
5. Defining the heap transition system allows us to reduce the problem to a reachability algorithm.

Prerequisite: Heap **List Fitness** algorithm (here *best fit* in ML-style syntax)

```
1: let best (cur:Chunk)(sz:int)(cand:Chunk) =
2:   if (cur.size ≤ sz and cur.sz - sz ≤ cand.sz - sz)
3:   then cur else cand;;

4: let rec findfit (choice: a → b → c → d)(l:list)(sz:int)(cand:Chunk) in
5:   match l with
6:   | [] → cand
7:   | [cur::tail] → (findfit tail sz (choice cur size cand));;

8: let rec FIT Lists sz = match Lists with
9:   | [] → ⊥
10:  | [cur::tail] → let res = (findfit best cur sz ⊥) in
11:                   match res with
12:                   | ⊥ → (fit tail sz)
13:                   | cur;;
```

The FRACSE calculus (part 1)

$$\frac{\text{size}(p) = x \quad \mathbf{FIT}(\mathcal{H}.\mathcal{L}, x) = l_1}{\mathbf{FREE}(p)}$$

$$\Gamma'_a[l_1] \leftarrow \Gamma_a[l_1] - 1 \quad \Gamma'_f[l_1] \leftarrow \Gamma_f[l_1] + 1$$

$$\frac{\mathbf{FIT}(\mathcal{H}.\mathcal{L}, x) = l_1}{p = \mathbf{ALLOC}(x)}$$

$$\Gamma'_a[l_1] \leftarrow \Gamma_a[l_1] + 1 \quad \Gamma'_f[l_1] \leftarrow \Gamma_f[l_1] - 1$$

$$\frac{\text{size}(p) = x \quad \mathbf{FIT}(\mathcal{H}.\mathcal{L}, x) = l_1 \quad \mathbf{FIT}(\mathcal{H}.\mathcal{L}, x + e) = l_2}{p_2 = \mathbf{REALLOC}(p_1, x + e)}$$

$$\Gamma'_a[l_1] \leftarrow \Gamma_a[l_1] - 1 \quad \Gamma'_f[l_1] \leftarrow \Gamma_f[l_1] + 1 \quad \Gamma'_a[l_2] \leftarrow \Gamma_a[l_2] + 1 \quad \Gamma'_f[l_2] \leftarrow \Gamma_f[l_2] - 1$$

The FRACSE calculus (part 2)

$$\begin{array}{c}
 \frac{\text{size}(p_1) = x_1 \quad \text{size}(p_2) = x_2 \quad \text{FIT}(\mathcal{H}.\mathcal{L}, x_1) = l_1 \quad \text{FIT}(\mathcal{H}.\mathcal{L}, x_2) = l_2 \quad \text{FIT}(\mathcal{H}.\mathcal{L}, x_3) = l_3}{p_3 = \mathbf{COALESCE}(p_1, p_2)} \\
 \frac{\Gamma'_f[l_1] \leftarrow \Gamma_f[l_1] - 1 \quad \Gamma'_f[l_2] \leftarrow \Gamma_f[l_2] - 1 \quad \Gamma'_f[l_3] = \Gamma_f[l_3] + 1}{\text{size}(p) = x \quad \text{FIT}(\mathcal{H}.\mathcal{L}, x) = l_1 \quad \text{FIT}(\mathcal{H}.\mathcal{L}, x - o) = l_2 \quad \text{FIT}(\mathcal{H}.\mathcal{L}, o) = l_3} \\
 \frac{\text{ADJ}(p_1, p_2) \quad \Gamma'_f[l_1] \leftarrow \Gamma_f[l_1] - 1 \quad \Gamma'_f[l_2] \leftarrow \Gamma_f[l_2] + 1 \quad \Gamma'_f[l_3] \leftarrow \Gamma_f[l_3] + 1}{\text{ADJ}(Top, p) \quad \Gamma'_f[l] \leftarrow \Gamma_f[l] + 1 \quad Top \leftarrow p} \\
 \frac{\text{FIT}(\mathcal{H}.\mathcal{L}, x) = l}{p = \mathbf{EXTEND}(x)}
 \end{array}$$

Pitfalls

- ▶ There can be multiple heaps (ex: one per thread). Heap selection is not defined in the FRACSE semantics. As **FIT** uses a heap parameter, it can handle multiple heaps easily.
- ▶ There can be multiple allocators within a process (ex: Windows front-end / back-end) driven by an **activation heuristic** for each bucket size. Adding such activation heuristic is a reasonable extension.
- ▶ FRACSE uses *lists*, some allocators use *arrays* (ex: JEMalloc)
- ▶ Heap meta-data is **abstracted** by design. Some exploit techniques still rely on meta-data corruption. We argue that due to internal checks in allocators, heap meta-data corruption as an exploit technique is dying.
- ▶ Non-deterministic heap behavior is not covered (ex: Die Hard allocator randomization, LFH subsegment randomization, etc). We need a probabilistic semantics to define this.
- ▶ This presentation only covers user-land heap allocators, no kernel heap allocator.

Summing up

- ▶ This work may be the first attempt at defining the formal semantics of heap allocators.
- ▶ Heap allocator implementations are so different that making generic heap analysis is a challenge.
- ▶ However, we can distinguish some common functionalities (split/coallesce/extend operations, list-based abstraction, heap selection, etc).
- ▶ Focusing on targeting user data and using a heap layout abstraction seems like the only generic way of exploiting the heap.
- ▶ FRACSE implementation is still going on. Its calculus may evolve based on experiments.

Thanks for attending!

Questions?

Mail: julien.vanegue@gmail.com

Twitter: @jvanegue

(Some) Related work

1. Smashing C++ VPTRS (Eric Landuyt)
2. VuDo Malloc tricks (Michel Kaempf)
3. Once upon a free (Scut)
4. Advanced DLMalloc Exploits (JP)
5. Malloc Maleficarum (Phantasmal Phantasmagoria)
6. The use of set_head to defeat the wilderness (g463)
7. Heap Feng Shui (Alex Sotirov)
8. Understanding the Low Fragmentation Heap (Chris Valasek)
9. The House Of Lore : PTmalloc exploitation (blackngel)
10. Pseudomonarchia Jemallocum (argp and huku)