



OWASP

AppSecResearch 2012

Heap Exploitation Abstraction by Example

Patroklos Argyroudis, Chariton Karamitas
{argp, huku}@census-labs.com

Census, Inc.

Who are we

- ✦ Patroklos Argyroudis, argp
 - ✦ Researcher at Census, Inc. (www.census-labs.com)
 - ✦ Topics: kernel/heap exploitation, auditing
- ✦ Chariton Karamitas, huku
 - ✦ Student at AUTh, intern at Census, Inc.
 - ✦ Topics: compilers, heap exploitation, maths

Outline

- ✦ Example: FreeBSD kernel memory allocator (UMA)
- ✦ Example: Linux kernel memory allocator (SLUB)
- ✦ Example: jemalloc userland memory allocator
- ✦ Abstracting heap exploitation

Related Work

- ✦ “Attacking the Core: Kernel Exploiting Notes” [1]
 - ✦ twiz, sgrakkyu, Phrack, 2007
 - ✦ Linux (heap), Solaris (stack)
- ✦ “Kernel Wars” [2]
 - ✦ signedness.org, Black Hat EU, 2007
 - ✦ *BSD (mbuf), Windows (stack)

Related Work

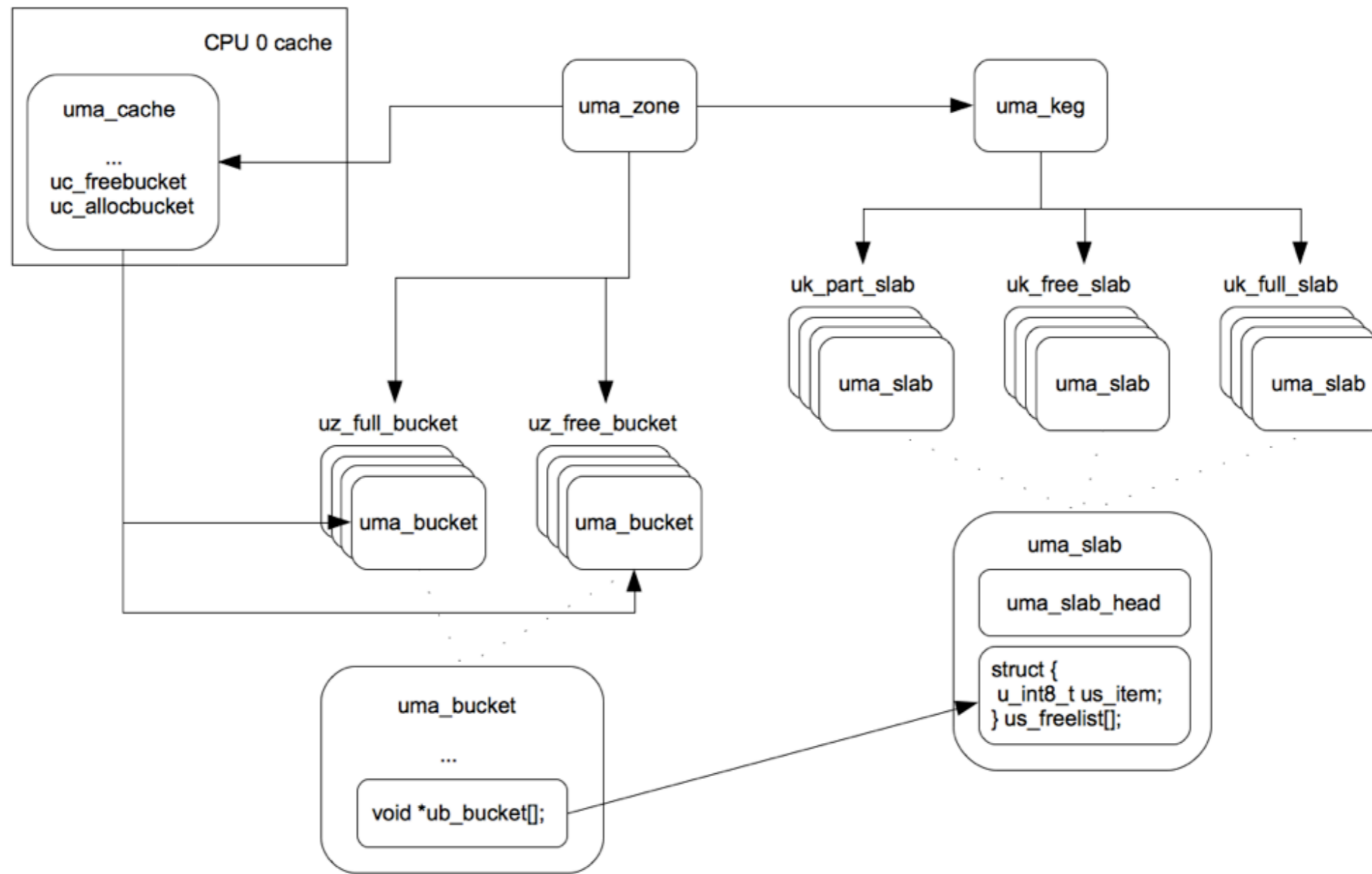
- ✦ “Exploitation in the Modern Era (Blueprint)” [3]
 - ✦ Chris Valasek, Ryan Smith, Black Hat EU, 2011
 - ✦ First attempt to abstract exploitation
- ✦ “Patras Heap Massacre” [4]
 - ✦ Chariton Karamitas, Patroklos Argyroudis, Fosscomm, 2011
 - ✦ Attempt to abstract heap exploitation

Example: FreeBSD UMA

Universal Memory Allocator

- ✦ FreeBSD's kernel memory allocator
 - ✦ Funded by Nokia for a proprietary project
 - ✦ The IPSO firewall/security appliance (thanks FX!)
 - ✦ Donated to FreeBSD
- ✦ Functions like a traditional slab allocator
 - ✦ Large areas, or slabs, of memory are pre-allocated
 - ✦ `malloc(9)` returns a free slot

UMA Architecture



UMA Architecture

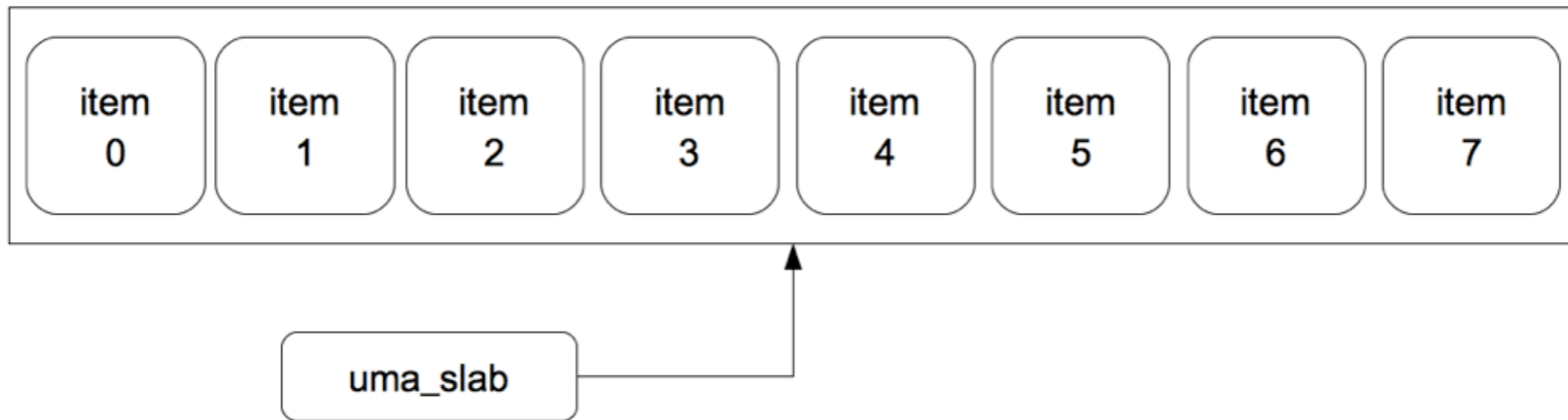
- ✦ Each zone (`uma_zone`) holds buckets (`uma_bucket`) of items
- ✦ The items are allocated on the zone's slabs (`uma_slab`)
- ✦ Each zone is associated with a keg (`uma_keg`)
- ✦ The keg holds the corresponding zone's slabs
- ✦ Each slab is of the same size as a page frame (usually 4096 bytes)
- ✦ Each slab has a slab header structure (`uma_slab_head`) which contains management metadata

vmstat(8)

ITEM	SIZE	LIMIT	USED	FREE	REQ	FAIL	SLEEP
UMA Kegs:	208,	0,	84,	1,	84,	0,	0
UMA Zones:	512,	0,	84,	0,	84,	0,	0
UMA Slabs:	568,	0,	872,	3,	3612,	0,	0
UMA RCntSlabs:	568,	0,	195,	1,	195,	0,	0
UMA Hash:	256,	0,	3,	12,	3,	0,	0
16 Bucket:	152,	0,	39,	11,	39,	0,	0
32 Bucket:	280,	0,	23,	5,	23,	0,	0
64 Bucket:	536,	0,	15,	6,	15,	57,	0
128 Bucket:	1048,	0,	19,	2,	19,	569,	0
VM OBJECT:	216,	0,	846,	72,	15779,	0,	0
MAP:	232,	0,	7,	25,	7,	0,	0
KMAP ENTRY:	120,	15004,	33,	122,	5952,	0,	0
MAP ENTRY:	120,	0,	543,	108,	35870,	0,	0
fakepg:	120,	0,	0,	0,	0,	0,	0
mt_zone:	4112,	0,	301,	10,	301,	0,	0
16:	16,	0,	2485,	203,	38918,	0,	0
32:	32,	0,	2780,	250,	19562,	0,	0
64:	64,	0,	5223,	153,	34106,	0,	0
128:	128,	0,	7156,	123,	10268,	0,	0
256:	256,	0,	830,	40,	3359,	0,	0
512:	512,	0,	347,	52,	3699,	0,	0
1024:	1024,	0,	50,	142,	6246,	0,	0
:█							

Slabs

An offpage slab of the "512" zone



A non-offpage slab of the "256" zone



uma_slab_head

```
238 struct uma_slab_head {
239     uma_keg_t      us_keg;           /* Keg we live in */
240     union {
241         LIST_ENTRY(uma_slab)  _us_link;    /* slabs in zone */
242         unsigned long  _us_size;    /* Size of allocation */
243     } us_type;
244     SLIST_ENTRY(uma_slab)  us_hlink;    /* Link for hash table */
245     u_int8_t      *us_data;           /* First item */
246     u_int8_t      us_flags;          /* Page flags see uma.h */
247     u_int8_t      us_freecount;    /* How many are free? */
248     u_int8_t      us_firstfree;    /* First free item index */
249 };
```


uma_keg

```
199 struct uma_keg {
200     LIST_ENTRY(uma_keg)    uk_link;        /* List of all kegs */
201
202     struct mtx             uk_lock;        /* Lock for the keg */
203     struct uma_hash       uk_hash;
204
205     char                   *uk_name;      /* Name of creating zone. */
206     LIST_HEAD(, uma_zone) uk_zones;      /* Keg's zones */
207     LIST_HEAD(, uma_slab) uk_part_slab;  /* partially allocated slabs */
208     LIST_HEAD(, uma_slab) uk_free_slab;  /* empty slab list */
209     LIST_HEAD(, uma_slab) uk_full_slab;  /* full slabs */
210
211     u_int32_t              uk_recurse;    /* Allocation recursion count */
212     u_int32_t              uk_align;      /* Alignment mask */
213     u_int32_t              uk_pages;      /* Total page count */
214     u_int32_t              uk_free;       /* Count of items free in slabs */
215     u_int32_t              uk_size;       /* Requested size of each item */
216     u_int32_t              uk_rsize;     /* Real size of each item */
217     u_int32_t              uk_maxpages;   /* Maximum number of pages to alloc */
218
219     uma_init               uk_init;      /* Keg's init routine */
220     uma_fini               uk_fini;      /* Keg's fini routine */
221     uma_alloc              uk_allocf;    /* Allocation function */
222     uma_free               uk_freef;     /* Free routine */
223
224     struct vm_object       *uk_obj;      /* Zone specific object */
225     vm_offset_t            uk_kva;       /* Base kva for zones with objs */
226     uma_zone_t             uk_slabzone;   /* Slab zone backing us, if OFFPAGE */
227
228     u_int16_t              uk_pgoff;     /* Offset to uma_slab struct */
229     u_int16_t              uk_ppera;     /* pages per allocation from backend */
230     u_int16_t              uk_ipers;     /* Items per slab */
231     u_int32_t              uk_flags;     /* Internal flags */
232 };
```

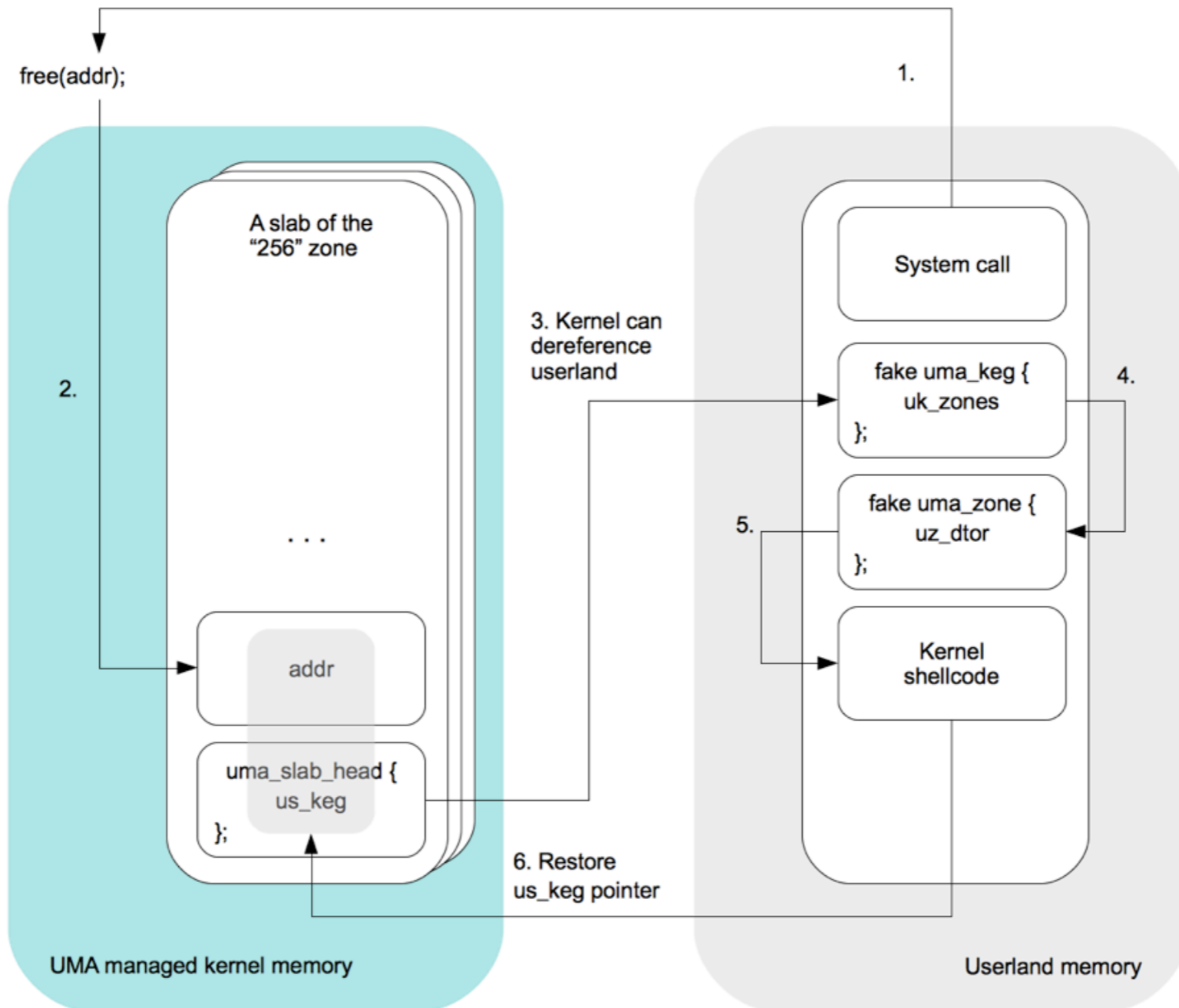

uma_zone

```
307 struct uma_zone {
308     char          *uz_name;          /* Text name of the zone */
309     struct mtx    *uz_lock;         /* Lock for the zone (keg's lock) */
310
311     LIST_ENTRY(uma_zone)  uz_link;    /* List of all zones in keg */
312     LIST_HEAD(, uma_bucket) uz_full_bucket; /* full buckets */
313     LIST_HEAD(, uma_bucket) uz_free_bucket; /* Buckets for frees */
314
315     LIST_HEAD(, uma_klink) uz_kegs;    /* List of kegs. */
316     struct uma_klink      uz_klink;    /* klink for first keg. */
317
318     uma_slaballoc  uz_slab;          /* Allocate a slab from the backend. */
319     uma_ctor       uz_ctor;         /* Constructor for each allocation */
320     uma_dtor       uz_dtor;         /* Destructor */
321     uma_init       uz_init;         /* Initializer for each item */
322     uma_fini       uz_fini;         /* Discards memory */
323
324     u_int32_t      uz_flags;        /* Flags inherited from kegs */
325     u_int32_t      uz_size;         /* Size inherited from kegs */
326
327     u_int64_t      uz_allocs UMA_ALIGN; /* Total number of allocations */
328     u_int64_t      uz_frees;        /* Total number of frees */
329     u_int64_t      uz_fails;        /* Total number of alloc failures */
330     u_int64_t      uz_sleeps;       /* Total number of alloc sleeps */
331     uint16_t       uz_fills;        /* Outstanding bucket fills */
332     uint16_t       uz_count;        /* Highest value ub_ptr can have */
333
334     /*
335     * This HAS to be the last item because we adjust the zone size
336     * based on NCPU and then allocate the space for the zones.
337     */
338     struct uma_cache uz_cpu[1]; /* Per cpu caches */
339 };
```


Code Execution

```
2528 uma_zfree_arg(uma_zone_t zone, void *item, void *udata)
2529 {
2530     uma_cache_t cache;
2531     uma_bucket_t bucket;
2532     int bflags;
2533     int cpu;
2534
2535 #ifdef UMA_DEBUG_ALLOC_1
2536     printf("Freeing item %p to %s(%p)\n", item, zone->uz_name, zone);
2537 #endif
2538     CTR2(KTR_UMA, "uma_zfree_arg thread %x zone %s", curthread,
2539         zone->uz_name);
2540
2541     /* uma_zfree(..., NULL) does nothing, to match free(9). */
2542     if (item == NULL)
2543         return;
2544
2545     if (zone->uz_dtor)
2546         zone->uz_dtor(item, zone->uz_size, udata);
```


uz_dtor Hijacking



Example: Linux SLUB

SLUB

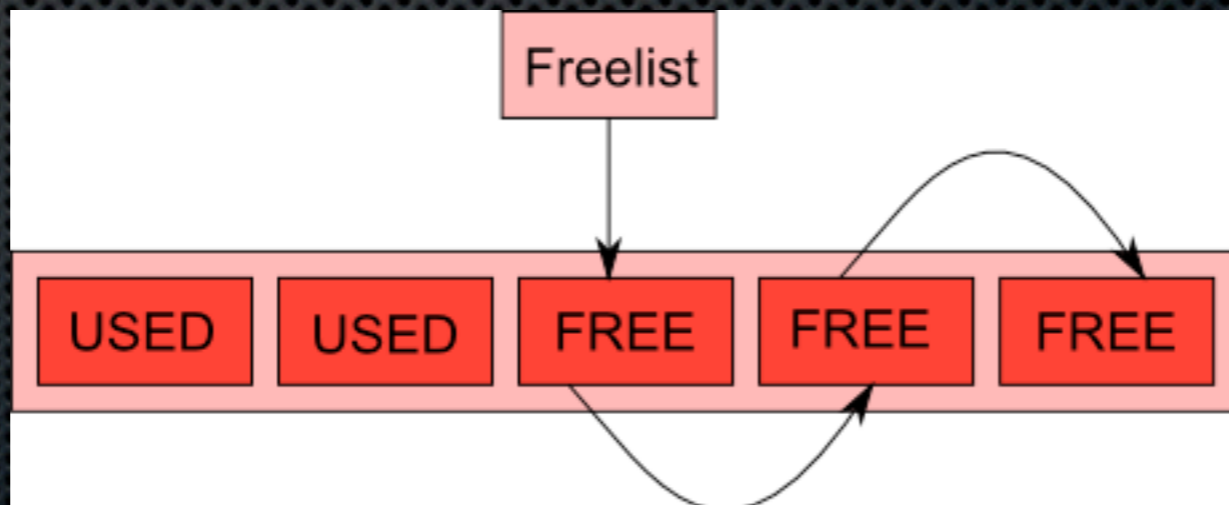
- ✦ Organizes physical memory frames in “caches” (UMA: kegs)
- ✦ Each cache holds slabs (UMA: slab) of objects (UMA: items) of the same size
 - ✦ kmalloc-32, kmalloc-64, task_struct, mm_struct
- ✦ Objects on a slab are contiguous
- ✦ A slab may have both allocated (used) and deallocated (free) objects

SLUB's slabs

- Each slab is at least PAGE_SIZE bytes (default 4096 bytes)
- A slab may span many pages
 - kmalloc-32: 128 objects * 32 bytes == 4096 bytes
 - task_struct (1088 bytes): 30 objects * 1088 bytes == 32640
 - A task_struct slab spans 8 pages
- Each CPU core has its own slabs

Metadata?

- No separate/dedicated metadata structures stored on the slabs
- Each free object stored on a slab has a next-free-object pointer
- Each slab has a page structure (struct page) that has a pointer (freelist) to the slab's first free object



SLUB's behavior

- ✦ Partial slabs: some free and some used objects
- ✦ New requests satisfied from partial slabs
 - ✦ Least-recently-used (LRU) policy
 - ✦ No partial slabs → allocation of new slab
- ✦ Generic slabs (e.g. kmalloc-32) are used to store different objects of the same size
 - ✦ Different kernel structures, buffers, etc
 - ✦ Contiguous

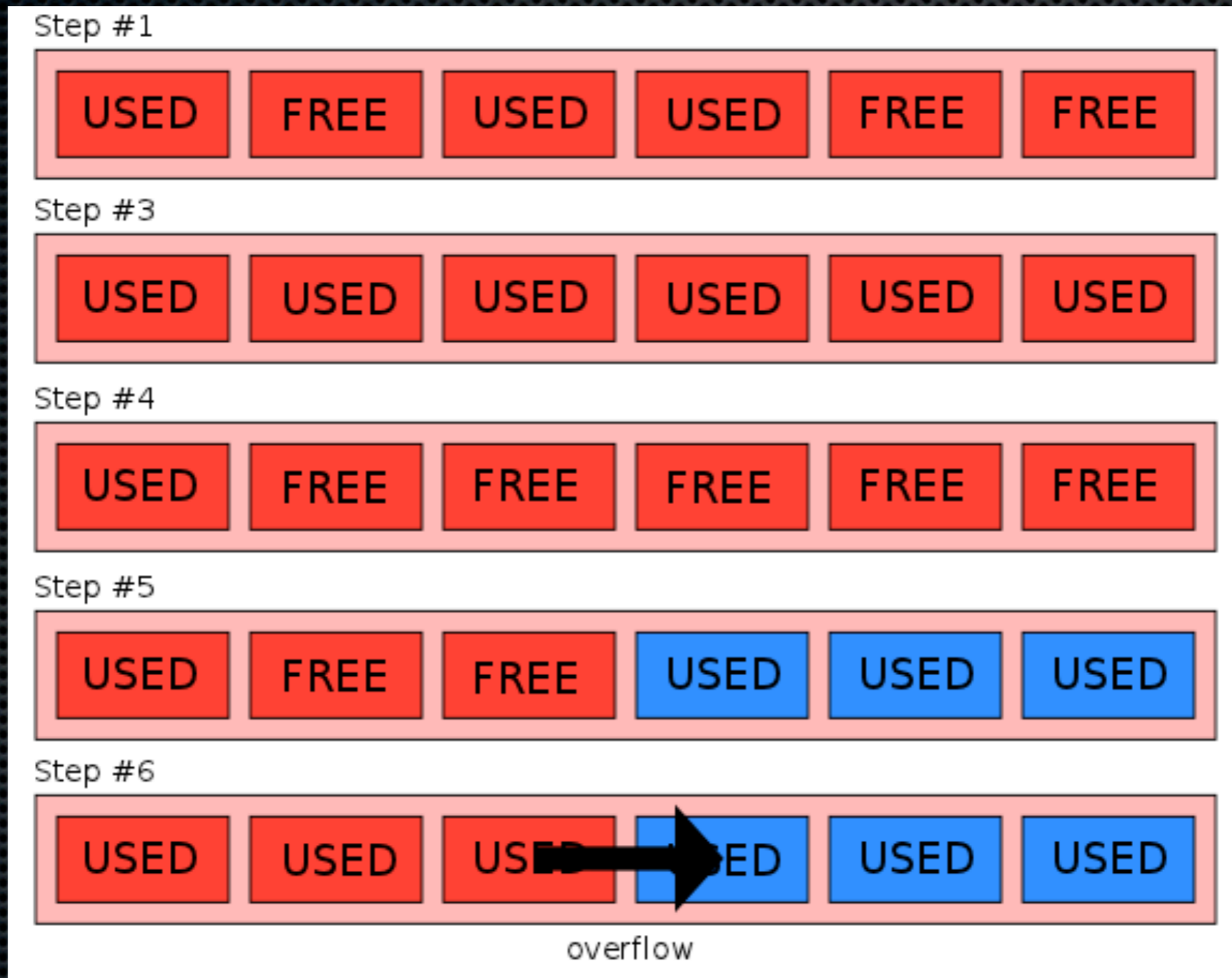
SLUB Exploitation

- ✦ Attack alternatives
 - ✦ Corrupt metadata of free objects on a slab
 - ✦ Corrupt adjacent objects on a slab
- ✦ We need a suitable kernel structure to corrupt
- ✦ We can allocate/deallocate from userland
- ✦ Same size as the object/structure we can overflow from
- ✦ Bring target slab to a predictable state in order to have the victim structure after the structure we can overflow from

SLUB Exploitation Algorithm

- Find free objects on target slab:
 - `cat /proc/slabinfo`
- Ensure allocations/deallocation happen on the slabs of the same CPU: `sched_setaffinity(2)`
- Consume a large number of objects that go on the target slab (reducing fragmentation)
- Deallocate a small number of objects from the target slab
- Allocate a smaller number of our selected victim objects
- Trigger the heap overflow bug overflowing onto the victim object

SLUB Exploitation



Victim Structure

- ✦ Traditionally struct `shmid_kernel`
- ✦ Allocations/deallocations controlled from userland
 - ✦ Allocation: `shmget(2)`
 - ✦ Deallocation: `ipcrm(1)`
- ✦ Leads to structure with yummy function pointers

shmid_kernel

```
86 struct shmid_kernel /* private to the kernel */
87 {
88     struct kern_ipc_perm    shm_perm;
89     struct file *          shm_file;
90     unsigned long          shm_nattch;
91     unsigned long          shm_segsz;
92     time_t                 shm_atim;
93     time_t                 shm_dtim;
94     time_t                 shm_ctim;
95     pid_t                  shm_cprid;
96     pid_t                  shm_lprid;
97     struct user_struct     *mlock_user;
98 };
```


file

```
934 struct file {
935     /*
936     * fu_list becomes invalid after file_free is called and queued via
937     * fu_rcuhead for RCU freeing
938     */
939     union {
940         struct list_head      fu_list;
941         struct rcu_head       fu_rcuhead;
942     } f_u;
943     struct path              f_path;
944 #define f_dentry             f_path.dentry
945 #define f_vfsmnt             f_path.mnt
946     const struct file_operations *f_op;
947     spinlock_t              f_lock; /* f_ep_links, f_flags, no IRQ */
948 #ifdef CONFIG_SMP
949     int                     f_sb_list_cpu;
950 #endif
951     atomic_long_t           f_count;
952     unsigned int            f_flags;
953     fmode_t                 f_mode;
954     loff_t                  f_pos;
955     struct fown_struct       f_owner;
956     const struct cred        *f_cred;
957     struct file_ra_state     f_ra;
958
959     u64                     f_version;
960 #ifdef CONFIG_SECURITY
961     void                    *f_security;
962 #endif
963     /* needed for tty driver, and maybe others */
964     void                    *private_data;
965
966 #ifdef CONFIG_EPOLL
967     /* Used by fs/eventpoll.c to link all the hooks to this file */
968     struct list_head        f_ep_links;
969 #endif /* #ifdef CONFIG_EPOLL */
970     struct address_space     *f_mapping;
971 #ifdef CONFIG_DEBUG_WRITECOUNT
972     unsigned long           f_mnt_write_state;
973 #endif
974 };
```


file_operations

```
1538 struct file_operations {
1539     struct module *owner;
1540     loff_t (*llseek) (struct file *, loff_t, int);
1541     ssize_t (*read) (struct file *, char __user *, size_t, loff_t *);
1542     ssize_t (*write) (struct file *, const char __user *, size_t, loff_t *);
1543     ssize_t (*aio_read) (struct kiocb *, const struct iovec *, unsigned long, loff_t);
1544     ssize_t (*aio_write) (struct kiocb *, const struct iovec *, unsigned long, loff_t);
1545     int (*readdir) (struct file *, void *, filldir_t);
1546     unsigned int (*poll) (struct file *, struct poll_table_struct *);
1547     long (*unlocked_ioctl) (struct file *, unsigned int, unsigned long);
1548     long (*compat_ioctl) (struct file *, unsigned int, unsigned long);
1549     int (*mmap) (struct file *, struct vm_area_struct *);
1550     int (*open) (struct inode *, struct file *);
1551     int (*flush) (struct file *, fl_owner_t id);
1552     int (*release) (struct inode *, struct file *);
1553     int (*fsync) (struct file *, int datasync);
1554     int (*aio_fsync) (struct kiocb *, int datasync);
1555     int (*fasync) (int, struct file *, int);
1556     int (*lock) (struct file *, int, struct file_lock *);
1557     ssize_t (*sendpage) (struct file *, struct page *, int, size_t, loff_t *, int);
1558     unsigned long (*get_unmapped_area) (struct file *, unsigned long, unsigned long, unsigned long, unsigned long);
1559     int (*check_flags) (int);
1560     int (*flock) (struct file *, int, struct file_lock *);
1561     ssize_t (*splice_write) (struct pipe_inode_info *, struct file *, loff_t *, size_t, unsigned int);
1562     ssize_t (*splice_read) (struct file *, loff_t *, struct pipe_inode_info *, size_t, unsigned int);
1563     int (*setlease) (struct file *, long, struct file_lock **);
1564     long (*fallocate) (struct file *file, int mode, loff_t offset,
1565                       loff_t len);
1566 };
```


Example: jemalloc

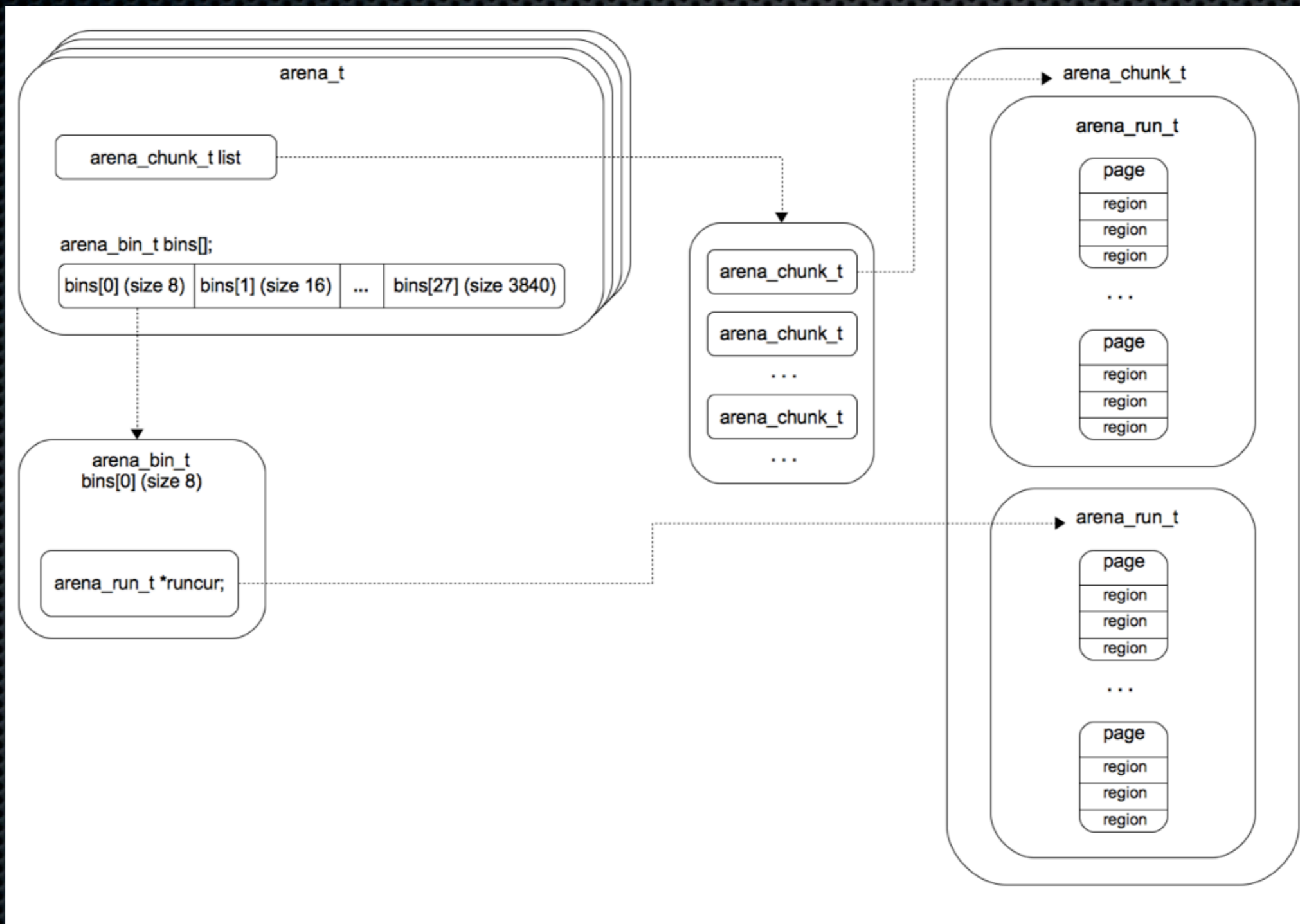
jemalloc

- ✦ FreeBSD needed a high performance, SMP-capable userland (libc) allocator
- ✦ Mozilla Firefox (Windows, Linux, Mac OS X)
- ✦ NetBSD libc
- ✦ Standalone version
- ✦ Facebook, to handle the load of its web services
- ✦ Defcon CTF is based on FreeBSD

jemalloc overview

- ✦ Memory is divided into chunks, always of the same size
- ✦ Chunks store all jemalloc data structures and user-requested memory (regions)
- ✦ Chunks are further divided into runs
- ✦ Runs keep track of free/used regions of specific sizes
- ✦ Regions are the heap items returned by malloc()
- ✦ Each run is associated with a bin, which stores trees of free regions (of its run)

jemalloc Architecture



jemalloc Exploitation

- ✦ Adjacent memory overwrite
- ✦ Metadata overwrite
 - ✦ Run header corruption
 - ✦ Chunk header corruption
 - ✦ Magazine (a.k.a thread cache) corruption
- ✦ For the details attend our Black Hat USA 2012 talk!

Abstracting Heap Exploitation

UMA - SLUB - jemalloc

- ✦ End-user allocations: UMA - items, SLUB - objects, jemalloc - regions
- ✦ Allocation containers: UMA - slabs, SLUB - slabs, jemalloc - runs
- ✦ Container groupings: UMA - kegs, SLUB - caches, jemalloc - chunks
- ✦ Execution-specific metadata:
 - ✦ UMA - zone, Linux kernel - zone, jemalloc - arena
 - ✦ UMA - buckets, SLUB - N/A, jemalloc - bins

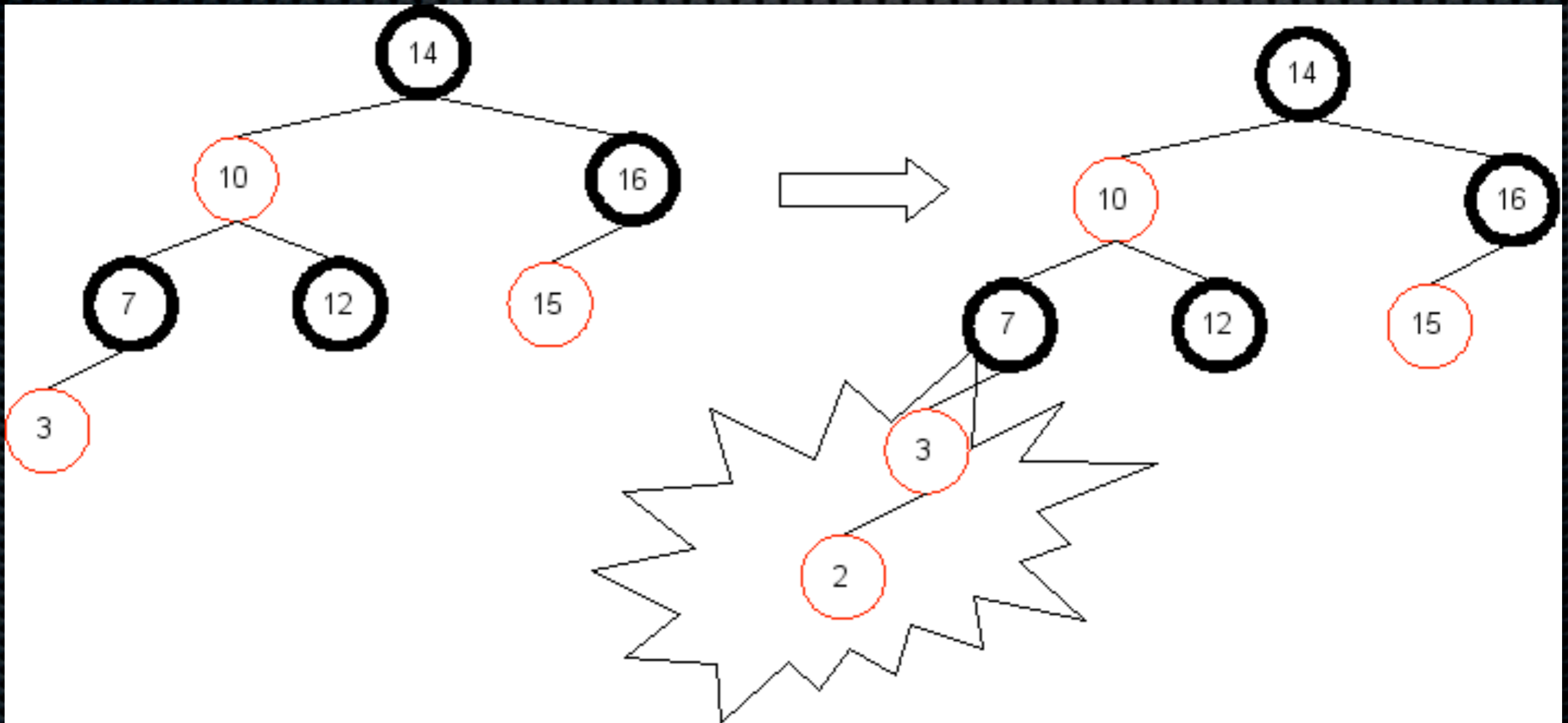
Value of Abstraction

- ✦ Chris Valasek's and Ryan Smith's Black Hat EU 2011 talk on abstracting exploitation through primitives [3]
- ✦ Back in CS 101 we were taught that abstraction is the most important skill of a computer scientist
- ✦ Specific exploitation techniques will become obsolete
- ✦ Our 2 drachmas are to abstract heap exploitation and have “primitives” that can be applied to new targets

Memory Allocators as Weird Machines

- ✦ Weird machine: The state machine of the target program after memory corruption [5, 6]
- ✦ In our case
 - ✦ State machine: Memory allocator
 - ✦ Weird machine: Post-corruption memory allocator
 - ✦ New states, unexpected by the developer
 - ✦ However reachable due to the memory corruption


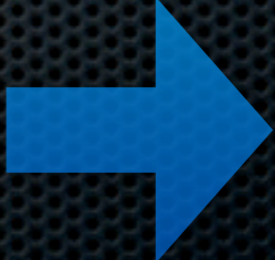
Heap Weird Machines



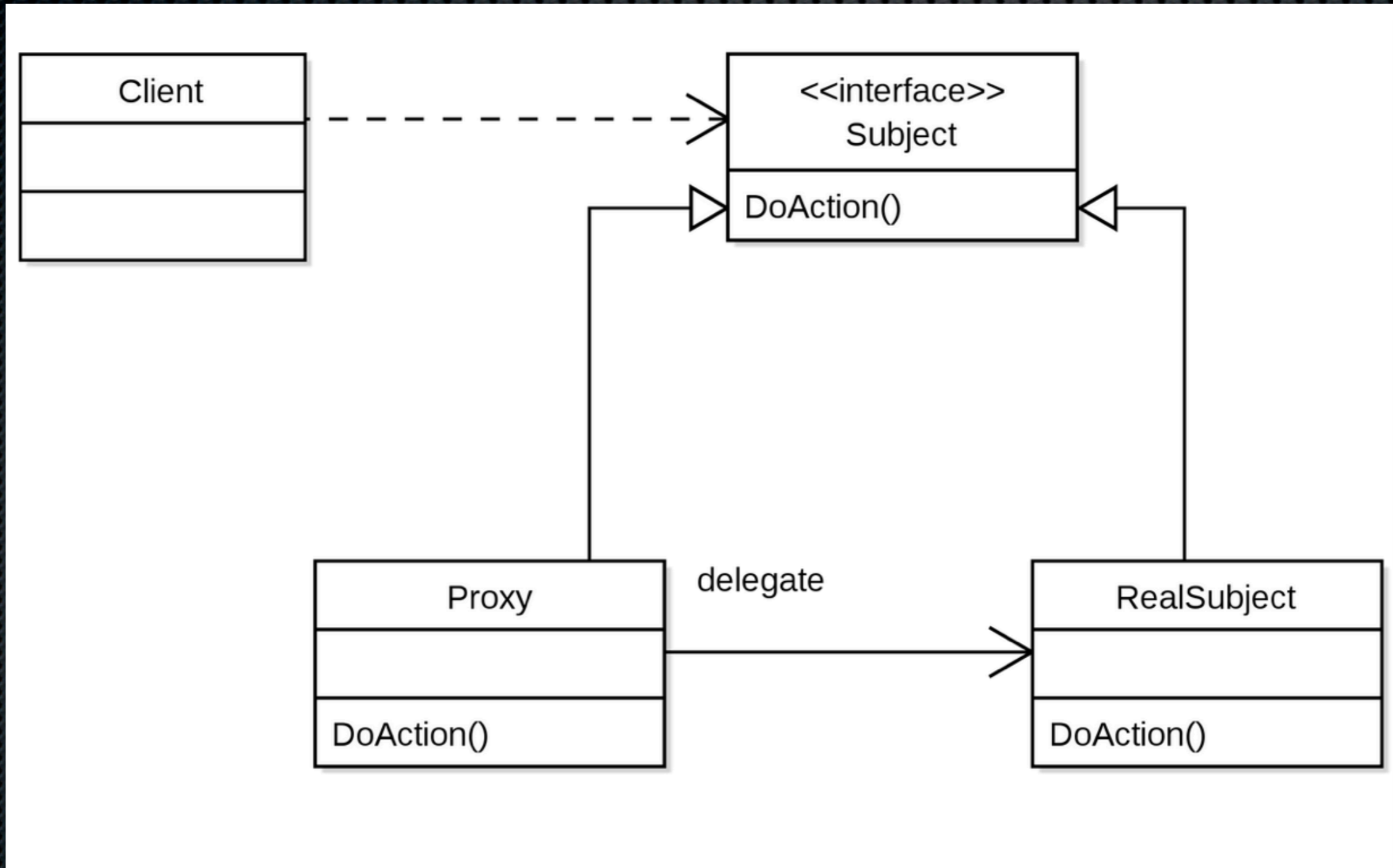
Heap Weird Machines

- ✦ Our memory allocator model: deterministic automaton (threads not taken into account)
- ✦ Metadata corruption abstraction
 - ✦ Corruption of the automaton's transition function
 - ✦ New states are reachable - most dead but not all
- ✦ Data (e.g. adjacent item) corruption abstraction
 - ✦ Manipulation of the automaton's determinacy
 - ✦ We control the order of transitions

The Weirding Module ;)

- ✦ The target heap manager should be treated as a high level API
 - ✦ For allocations and deallocations
- ✦ “Applications” that use the allocator (Javascript, system calls, incoming packets) provide a way to proxy these API calls
- ✦ Attacker  Application (Proxy)  Allocator

The Weirding Module ;)



Conclusion

- ✦ Future work
 - ✦ Operational semantics (formal notation)
 - ✦ More examples on both allocators and exploits
- ✦ Acknowledgments
 - ✦ Dr ;) Dimitris Glynos
 - ✦ Chris Valasek
 - ✦ Sergey Bratus

References

- [1] “Attacking the Core: Kernel Exploiting Notes”, twiz, sgrakkyu, Phrack, 2007
- [2] “Kernel Wars”, signedness.org, Black Hat EU, 2007
- [3] “Exploitation in the Modern Era (Blueprint)”, Chris Valasek, Ryan Smith, Black Hat EU, 2011
- [4] “Patras Heap Massacre”, Chariton Karamitas, Patroklos Argyroudis, Fosscomm, 2011
- [5] “Exploit Programming”, Sergey Bratus et al, ;login:, 2011
- [6] “Exploitation and State Machines”, Halvar Flake, Infiltrate, 2011