

# CS380L: Advanced Operating Systems Lab #2

Zeyuan Hu <sup>1</sup>, iamzeyuanhu@utexas.edu

EID:zh4378 Spring 2019

## 1 Environment

We have two Linux machines for experiments: `erode-lom` and `thoothukudi-lom`. Both of them run Ubuntu 16.04 LTS. <sup>2</sup>

## 2 System Tools Exercise

### 2.1 strace

`script` command allows user to record terminal printout into a file [1]. Per the lab instruction, we use `strace` to trace the syscalls and signals of a target process [2]. In our case, we trace the process involving `cat`. One thing I notice is that `script` contains some unicode as shown in Figure ??<sup>1</sup>. Thus, we use the following code to clean up the output:

```
cat $FILE | perl -pe 's/\e([\[\]]|\.[a-zA-Z]|\.a)/g' | col -b > $FILE-processed
```

The result is shown in Figure 1 <sup>3</sup>.

### 2.2 lsof

`lsof` lists all open files [3]. `lsof | grep /dev` shows all the open devices used by user-space programs. On our machine, we have the following opened devices:

- `/dev/null`: null device
- `/dev/pts/*` and `/dev/tty`: terminal devices
- `/dev/urandom`: kernel random number source device
- `/dev/ptmx`: a character file to create a pseudoterminal master

---

<sup>1</sup>30 hours spent on this lab.

<sup>2</sup>`erode-lom` with ip address: 192.168.1.120 and `thoothukudi-lom` with ip address: 192.168.1.121. Due to the network setup, those two machines are not properly registered in the department DNS server. We use ip address instead host names to logon these two machines.

<sup>3</sup>raw output and cleanup output comes with the report as `session.record` and `session.record-processed` respectively

```

Script started on Fri 01 Mar 2019 04:37:34 PM CST
ESC[1;34m zeyuanhu @ ESC[0;36mHotDog(thoothukudi-lom)ESC[0mESC[1;34m ~ESC[0m
ESC[0;36m Fri Mar 01 16:37:34 $ ESC[0;39mstrace cat - > new_file
execve("/bin/cat", ["cat", "-"], [/ * 37 vars */]) = 0
brk(NULL)                                = 0x2280000
access("/etc/ld.so.nohwcap", F_OK)        = -1 ENOENT (No such file or directory)
access("/etc/ld.so.preload", R_OK)        = -1 ENOENT (No such file or directory)

```

Raw output of script, which contains Unicode character.

```

Script started on Fri 01 Mar 2019 04:37:34 PM CST
zeyuanhu @ HotDog(thoothukudi-lom) ~
Fri Mar 01 16:37:34 $ strace cat - > new_file
execve("/bin/cat", ["cat", "-"], [/ * 37 vars */]) = 0
brk(NULL)                                = 0x2280000
access("/etc/ld.so.nohwcap", F_OK)        = -1 ENOENT (No such file or directory)
access("/etc/ld.so.preload", R_OK)        = -1 ENOENT (No such file or directory)

```

Script output after cleanup

Figure 1: script output before and after cleanup

### 3 Network Tools

`ifconfig` command lists all the network interfaces the machine is using to communicate externally. On our machine, interface for Ethernet is `eno1`. We can find IP address, gateway address, and subnet mask from the output.

`tcpdump` command can dump traffic on a network interface. We follow the instruction and use two packets to capture packets. On one terminal, we run `sudo tcpdump -i eno1 -w tcpdump.out.1` and on the other terminal, we first release our machine's current DHCP lease by executing `dhclient -r` and recontact DHCP server via `sudo dhclient eno1`. We stop capturing once `sudo dhclient eno1` finishes. We answer the following questions based on the packets we record.

#### a. Are DHCP messages sent over UDP or TCP?

We use `tcpdump -nn -r tcpdump.out.1 | grep -i dhcp` to filter out the DHCP messages from the dump. `-nn` ensures that we can see the actual port number instead of the port name. One line of the output is:

```

19:07:21.409099 IP 0.0.0.0.68 > 255.255.255.255.67: BOOTP/DHCP, Request from
10:98:36:ac:9d:45, length 300

```

The first field shows the time that the packet was traveling. The second field shows the source host address and port, followed by the destination host address and port. The third field shows the protocol the packet was using. From DHCP [4], we know DHCP messages sent over UDP. As shown by the printout, the messages are sent between port 68 (client) and port 67 (server).

- b. What is the link-layer (e.g., Ethernet) address of your host? (Feel free to obscure the last couple bytes for privacy's sake)**

We use the same `tcpdump` command as above with extra `-e` option to show link-layer header. The following printout contains DHCP messages for acquiring IP address:

```
19:07:41.375989 10:98:36:ac:9d:45 > ff:ff:ff:ff:ff:ff, ethertype IPv4 (0x0800
    ), length 342: 0.0.0.0.68 > 255.255.255.255.67: BOOTP/DHCP, Request from
    10:98:36:ac:9d:45, length 300
19:07:44.654365 d4:ae:52:e8:4e:85 > 10:98:36:ac:9d:47, ethertype IPv4 (0x0800
    ), length 342: 192.168.1.1.67 > 192.168.1.221.68: BOOTP/DHCP, Reply,
    length 300
```

From the printout we can see that the link-layer address (MAC address) of the host is 10:98:36:ac:9d:45.

- c. What is the IP address of your DHCP server?**

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

- d. What is the purpose of the DHCP release message? What would happen if the client's DHCP release message is lost?**

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

- e. Does the DHCP server issue an acknowledgment of receipt of the client's DHCP request?**

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum

L1-dcache-load-misses	[Hardware cache event]
L1-dcache-loads	[Hardware cache event]
L1-dcache-stores	[Hardware cache event]
L1-icache-load-misses	[Hardware cache event]
LLC-load-misses	[Hardware cache event]
LLC-loads	[Hardware cache event]
LLC-store-misses	[Hardware cache event]
LLC-stores	[Hardware cache event]
branch-load-misses	[Hardware cache event]
branch-loads	[Hardware cache event]
dTLB-load-misses	[Hardware cache event]
dTLB-loads	[Hardware cache event]
dTLB-store-misses	[Hardware cache event]
dTLB-stores	[Hardware cache event]
iTLB-load-misses	[Hardware cache event]
iTLB-loads	[Hardware cache event]
node-load-misses	[Hardware cache event]
node-loads	[Hardware cache event]
node-store-misses	[Hardware cache event]
node-stores	[Hardware cache event]

Figure 2: Part of output of `perf list` related to cache

sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

**f. What would happen if the client’s DHCP release message is lost?**

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

## 4 perf\_event\_open

We first check the support of `perf_event_open` interface by checking the existence of `/proc/sys/kernel/perf_event_paranoid`, which is true in our case (value is set to 2). In Linux, `perf_event_open` interface [?] is used to setup performance monitoring. Specifically, it provides an interface that allows user to access various events (i.e., events counted by performance counters [5]). `perf list` gives available events on current machine. Counters related to cache in our machine is shown in Figure 2. We are interested in counters related to L1 data cache and data TLB. As shown in Figure 2, we have counters for number of times the L1 cache was accessed for data (`L1-dcache-loads`), the number of those access that resulted in a cache miss (`L1-dcache-load-misses`), and write access of L1 data cache (`L1-dcache-stores`). Similarly, for data TLB, we have read access (`dTLB-loads`), read miss (`dTLB-load-misses`), write access (`dTLB-stores`), and write miss (`dTLB-store-misses`).

```

000000000400d36 <perf_event_open>:
400d36: 55                push   %rbp
400d37: 48 89 e5          mov    %rsp,%rbp
400d3a: 48 83 ec 20       sub    $0x20,%rsp
400d3e: 48 89 7d f8       mov    %rdi,-0x8(%rbp)
400d42: 89 75 f4         mov    %esi,-0xc(%rbp)
400d45: 89 55 f0         mov    %edx,-0x10(%rbp)
400d48: 89 4d ec         mov    %ecx,-0x14(%rbp)
400d4b: 4c 89 45 e0       mov    %r8,-0x20(%rbp)
400d4f: 48 8b 7d e0       mov    -0x20(%rbp),%rdi
400d53: 8b 75 ec         mov    -0x14(%rbp),%esi
400d56: 8b 4d f0         mov    -0x10(%rbp),%ecx
400d59: 8b 55 f4         mov    -0xc(%rbp),%edx
400d5c: 48 8b 45 f8       mov    -0x8(%rbp),%rax
400d60: 49 89 f9         mov    %rdi,%r9
400d63: 41 89 f0         mov    %esi,%r8d
400d66: 48 89 c6         mov    %rax,%rsi
400d69: bf 2a 01 00 00   mov    $0x12a,%edi
400d6e: b8 00 00 00 00   mov    $0x0,%eax
400d73: e8 c8 fd ff ff   callq 400b40 <syscall@plt>
400d78: c9              leaveq %eax
400d79: c3              retq

```

Figure 3: perf\_event\_open section of objdump -d output

According to man page [?], there is no glibc wrapper for perf\_event\_open system call. However, we can wrap it on our own as the following:

```

1 static long perf_event_open(struct perf_event_attr *hw_event, pid_t pid, int cpu,
   int group_fd, unsigned long flags) {
2     return syscall(__NR_perf_event_open, hw_event, pid, cpu, group_fd, flags);
3 }

```

Using the syscall does not mean there is a syscall opcode in the program. Figure 3 shows perf\_event\_open section of objdump -d. Using syscall function will have perf\_event\_open system call number compiled in the program. In this case, the system call number of perf\_event\_open is 298 (can be checked in arch/x86/entry/syscalls/syscall\_64.tbl of the Linux kernel source tree), which is \$0x12a in hex.

We want to monitor the following events: read, write, read miss for level 1 data cache and read miss, write miss for data TLB. A call to perf\_event\_open gives a file descriptor that corresponds to one event being measured. We can use ioctl interface to control the counters and read file descriptors to obtain counter vales. The following code highlights the usage of perf\_event\_open interface to measure all five events for trivial printf <sup>4</sup>.

```

1 int hw_cache_perf_event_open(int group_fd, int cache_id, int cache_op_id,
2 int cache_op_result_id) {
3     struct perf_event_attr pe;
4     memset(&pe, 0, sizeof(struct perf_event_attr));

```

<sup>4</sup>complete code can be seen in perf\_event\_open\_usage.c

```

5  pe.type = PERF_TYPE_HW_CACHE;
6  pe.size = sizeof(struct perf_event_attr);
7  pe.config = cache_id | (cache_op_id << 8) | (cache_op_result_id << 16);
8  pe.disabled = 0;
9  if (group_fd == -1) {
10     pe.disabled = 1;
11 }
12 pe.exclude_kernel = 1;
13 pe.exclude_hv = 1;
14 int fd = perf_event_open(&pe, 0, cpu_id, group_fd, 0);
15 if (fd == -1) {
16     perror("perf_event_open");
17     exit(EXIT_FAILURE);
18 }
19 return fd;
20 }
21
22 int main(int argc, char **argv) {
23
24     int l1_read_access_fd = hw_cache_perf_event_open(
25         -1, PERF_COUNT_HW_CACHE_L1D, PERF_COUNT_HW_CACHE_OP_READ,
26         PERF_COUNT_HW_CACHE_RESULT_ACCESS);
27     int leader_fd = l1_read_access_fd;
28
29     int l1_read_miss_fd = hw_cache_perf_event_open(
30         leader_fd, PERF_COUNT_HW_CACHE_L1D, PERF_COUNT_HW_CACHE_OP_READ,
31         PERF_COUNT_HW_CACHE_RESULT_MISS);
32     int l1_write_access_fd = hw_cache_perf_event_open(
33         leader_fd, PERF_COUNT_HW_CACHE_L1D, PERF_COUNT_HW_CACHE_OP_WRITE,
34         PERF_COUNT_HW_CACHE_RESULT_ACCESS);
35     int tlb_read_miss_fd = hw_cache_perf_event_open(
36         leader_fd, PERF_COUNT_HW_CACHE_DTLB, PERF_COUNT_HW_CACHE_OP_READ,
37         PERF_COUNT_HW_CACHE_RESULT_MISS);
38     int tlb_write_miss_fd = hw_cache_perf_event_open(
39         leader_fd, PERF_COUNT_HW_CACHE_DTLB, PERF_COUNT_HW_CACHE_OP_WRITE,
40         PERF_COUNT_HW_CACHE_RESULT_MISS);

```

```

41
42     ioctl(leader_fd, PERF_EVENT_IOC_RESET, PERF_IOC_FLAG_GROUP);
43     ioctl(leader_fd, PERF_EVENT_IOC_ENABLE, PERF_IOC_FLAG_GROUP);
44
45     // Do the work that we want to analyze
46     printf("Do some work that we want to measure here\n");
47
48     ioctl(leader_fd, PERF_EVENT_IOC_DISABLE, PERF_IOC_FLAG_GROUP);
49
50     uint64_t l1_read_miss = 0;
51     uint64_t l1_read_access = 0;
52     uint64_t l1_write_access = 0;
53     uint64_t tlb_read_miss = 0;
54     uint64_t tlb_write_miss = 0;
55
56     read(l1_read_access_fd, &l1_read_access, sizeof(uint64_t));
57     read(l1_read_miss_fd, &l1_read_miss, sizeof(uint64_t));
58     read(l1_write_access_fd, &l1_write_access, sizeof(uint64_t));
59     read(tlb_read_miss_fd, &tlb_read_miss, sizeof(uint64_t));
60     read(tlb_write_miss_fd, &tlb_write_miss, sizeof(uint64_t));
61
62     close(l1_read_access_fd);
63     close(l1_read_miss_fd);
64     close(l1_write_access_fd);
65     close(tlb_read_miss_fd);
66     close(tlb_write_miss_fd);
67
68     printf("[Performance counters]\n");
69     printf("Data L1 read access: %" PRIu64 "\n", l1_read_access);
70     printf("Data L1 write access: %" PRIu64 "\n", l1_write_access);
71     printf("Data L1 read miss: %" PRIu64 "\n", l1_read_miss);
72     printf("Data L1 read miss rate: %.5f\n",
73           (double)l1_read_miss / l1_read_access);
74     printf("Data TLB read miss: %" PRIu64 "\n", tlb_read_miss);
75     printf("Data TLB write miss: %" PRIu64 "\n", tlb_write_miss);
76

```

```

77     fflush(stdout);
78
79     return 0;
80 }

```

File descriptors returned from calling `perf_event_open` can be grouped together so that we can measure corresponding events simultaneously. There are five events we want to measure and as shown in the code above, we group them together and measure them at the same time. However, measuring all five events at the same time may not be possible for some machine as CPU only has limited amount of machine specific registers (MSRs) for low-level performance counting. In our environment, this number is 8<sup>5</sup>. For the following experiment, we also want to lock our process onto a single processor. To achieve it, we use the following code:

```

1  cpu_set_t set;
2  CPU_ZERO(&set);
3  CPU_SET(cpu_id, &set);
4  if (sched_setaffinity(0, sizeof(cpu_set_t), &set) == -1) {
5      perror("sched_setaffinity");
6      exit(EXIT_FAILURE);
7  }

```

Here, all counters we setup are associated with CPU specified by `cpu_id`. We use `sched_setaffinity` system call to set the CPU affinity of current process and ensure that it runs on CPU of `cpu_id` only.

## 5 Measuring memory access behavior

`mmap` system call maps files or devices into memory. The created mapping can be file-backed or anonymous. File-backed mapping maps an area of the process's virtual memory to files (i.e., an area of the process's virtual memory is mapped to file-backed memory); reading those areas of memory causes the file read. Anonymous mapping is the opposite of file-backed mapping (i.e., not backed by file; an area of the process's virtual memory is mapped to anonymous memory). In this section, we experiment with both types of mappings. Since file-backed mapping can be shared or private, we have three memory mappings setup in our experiment: anonymous, file-based (private) and file-based (share). In this experiment, we use `MAP_POPULATE` flag. `MAP_POPULATE` populates page tables for a mapping. For file-backed mapping, this means read-ahead on file, which reduce blocking on

---

<sup>5</sup>check via `cpuid|grep 'number of counters per logical processor'`



page faults later. In addition, we use `memset` to initialize the mapped region and `msync` to flush the change made to the file-backed memory back to file system whenever possible.

We study the following code, which is used to access the mapped region:

```
1 #define CACHE_LINE_SIZE 64
2 int opt_random_access;
3 void do_mem_access(char *p, int size) {
4     int outer, locality, i;
5     int ws_base = 0;
6     int max_base = size / CACHE_LINE_SIZE - 512;
7     for (outer = 0; outer < (1 << 20); outer++) {
8         long r = simplerand() % max_base;
9         if (opt_random_access) {
10             ws_base = r;
11         } else {
12             ws_base += 512;
13             if (ws_base >= max_base) {
14                 ws_base = 0;
15             }
16         }
17         for (locality = 0; locality < 16; locality++) {
18             volatile char *a;
19             char c;
20             for (i = 0; i < 512; i++) {
21                 a = p + ws_base + i * CACHE_LINE_SIZE;
22                 if (i % 8 == 0) {
23                     *a = 1;
24                 } else {
25                     c = *a;
26                 }
27             }
28         }
29     }
30 }
```

`simplerand()` is a function generates random number. A global variable `opt_random_access` is used to specify the memory access to be sequential or random. Specifically, `db_mem_access()`

	[sequential access]	[random access]
Data L1 read access:	7521439271.600 (std 227.930)	7521442554.200 (std 142.106)
Data L1 write access:	1077936131.000 (std 0.000)	1077936131.000 (std 0.000)
Data L1 read miss:	289466866.800 (std 198047.028)	826536782.800 (std 393229.625)
Data L1 read miss rate:	0.038 (std 0.000)	0.110 (std 0.000)
Data TLB read miss:	99193.400 (std 717.220)	1551880.800 (std 1346.083)
Data TLB write miss:	32087.400 (std 433.014)	737437.800 (std 736.102)
utime	3.317 (std 0.005) s	4.827 (std 0.003) s
stime	0.006 (std 0.005) s	0.008 (std 0.001) s
minflt	4113.600 (std 0.490)	7336.600 (std 0.490)
maxrss	16637.600 (std 42.828) KB	16744.000 (std 41.800) KB

Figure 4: The statistics of sequential and random memory accesses with anonymous mapped memory regions

routine will access a working set consists of 512 consecutive cache lines. Among those 512 cache lines,  $\frac{1}{8}$  will be write access and  $\frac{7}{8}$  will be read access. The same working set is accessed 16 times. The whole routine will work on  $2^{20}$  working sets. Working sets are selected either sequentially or in random (specified by `opt_random_access`). In sequential access case, offset between working sets is 512 (e.g., since there are 512 cache lines of a set). In random access case, the start(i.e., base) of working set is selected randomly based on the number generated by `simplerand()`. `simplerand()` is called no matter the sequential memory access or random access because we want to keep the experiment results difference only due to the memory access pattern difference. In other words, we want to have both memory access patterns have the same function call and context switch overheads (caused by calling `simplerand()`).

We study `db_mem_access()` behavior on a VM with 2GB of memory and 4 virtual CPU cores. The VM runs Ubuntu 18.10 with kernel version 4.19.6. `-cpu host` is added as part of QEMU option to ensure our VM can access performance counters. We allocate 1GB buffer via `mmap` as input for `do_mem_access`. To ensure the repeatable result, we flush the level 1 data cache by accessing a 16MB buffer before calling `do_mem_access`. We invoke `getrusage` just before and after `do_mem_access` to collect resource usage statistics of initialization code and `do_mem_access`. We run the program 5 times and calculate mean and stand deviation of statistics <sup>6</sup>. We also submit the raw results along with this report.

We first compare sequential access and random access (indicated by `opt_random_access`). We use anonymous mapped memory regions. Figure 4 shows statistics for both cases <sup>7</sup>. As shown in the figure, the ratio between read and write access is 7:1 for data L1 cache for both cases, which is as expected. The read miss rate in random access case is 3x higher than sequential access. For data TLB, read miss and write miss of random access is 15.6X and 23x higher than sequential access respectively. Random access runs 1.46x slower than sequential access. From those statistics, we can

<sup>6</sup>We compile our program with `-O2` to prevent any unwanted memory access from the generated executable

<sup>7</sup>`majflt`, `inblock`, `oublock` are omitted because they are all 0

	[private, seq]	[private, rand]	[share, seq]	[share, rand]
maxrss (init)	17540.800 (std 17.046) KB	17545.600 (std 18.693) KB	17550.400 (std 23.269) KB	17516.800 (std 15.052) KB
inblock (init)	0.000 (std 0.000)	0.000 (std 0.000)	0.000 (std 0.000)	0.000 (std 0.000)
oublock (init)	0.000 (std 0.000)	0.000 (std 0.000)	0.000 (std 0.000)	0.000 (std 0.000)
maxrss (fun)	16658.400 (std 45.789) KB	16816.000 (std 52.642) KB	16630.400 (std 45.438) KB	16839.200 (std 19.823) KB
inblock (fun)	0.000 (std 0.000)	0.000 (std 0.000)	0.000 (std 0.000)	0.000 (std 0.000)
oublock (fun)	0.000 (std 0.000)	0.000 (std 0.000)	6566.400 (std 13132.800)	6566.400 (std 13132.800)

Figure 5: The statistics of sequential and random memory accesses with file-backed mapping with MAP\_PRIVATE and MAP\_SHARED. fun is a shorthand for do\_mem\_access.

	[private, pop]	[share, pop]	[share, memset]
maxrss (init)	1066060.000 (std 106.883) KB	1064109.600 (std 274.140) KB	1066129.600 (std 31.455) KB
inblock (init)	509718.400 (std 224734.764)	80596.800 (std 161193.600)	0.000 (std 0.000)
oublock (init)	0.000 (std 0.000)	0.000 (std 0.000)	2064320.000 (std 0.000)
maxrss (fun)	0.000 (std 0.000)	2334.400 (std 259.200)	320.800 (std 34.260)
inblock (fun)	99.200 (std 198.400)	0.000 (std 0.000)	0.000 (std 0.000)
oublock (fun)	0.000 (std 0.000)	6566.400 (std 13132.800)	32832.000 (std 0.000)

Figure 6: The statistics of file-backed mapping for MAP\_POPULATE and memset after buffer allocation. fun is a shorthand for do\_mem\_access.

see that there is no significant difference in terms of runtime and miss rate of data L1 cache between random access and sequential access (e.g., miss rate for random access is around 0.1). In the case of sequential access, the number of data TLB read miss is about 99193, which is 396MB for 4KB page size. 396MB is smaller than the size of our buffer. If we take a look at the code, we notice `max_base` in `do_mem_access` restricts the maximum base address of the working set. If we take a look at how `max_base` is initialized and the line we initialize `a`, we can see that the working set is essentially the first  $\frac{1}{64}$  of the given buffer, which is 16MB, which matches with the value of `maxrss` we observed and is much lower than 1GB. Realizing this fact also explains why TLB miss count is lower than the total number of data pages <sup>8</sup>.

In the next step, we experiment with file-based memory mappings. In this part, we only concern `maxrss`, `inblock`, and `oublock` because the statistics related to data L1 cache and data TLB are generated with the same setup as anonymous mappings. In addition, statistics related to memory cache is not affected by the file-backed mapping. The statistics we have is shown in Figure 5. Compare to previous result, we notice there are large amount of `oublock` during `do_mem_access` for file-backed mapping with MAP\_SHARED. This is intuitive as the changes to the file have to be saved in shared mode. We notice that there is `inblock` is zero, which is due to the internal cache in Linux kernel for reading file. `oublock` for both random access and sequential access under shared mode are exactly the same. This may due to the fact that access pattern is not significant for the shared mode as synchronization for internal file cache is needed no matter what memory access pattern we may encounter.

In this experiment, we further examine MAP\_POPULATE flag and `memset` the mapped memory for initialization. Figure 6 shows the statistics. Note that all the columns are from the sequential

<sup>8</sup>If we want to access all the allocated buffer (1GB), we should replace line 21 (`a` initialization) with `a = p + (ws_base + i) * CACHE_LINE_SIZE;`

memory access. With `MAP_POPULATE` flag, the mapped memory will be populate page tables at the very beginning and perform read-ahead on the file. Thus we see the `maxrss` will match the allocated buffer size, which is around 1GB during the initialization. If the mapped memory is initialized with `memset` and then perform synchronization to the file system via `msync`, then there will be large amount of `oublock` during the initialization.

We also using `strace` to trace our program <sup>9</sup>. We find a system call `arch_prctl(ARCH_SET_FS, 0x7f83a006d540)` in the printout. The system call means setting FS register to `0x7f83a006d540` [6]. FS register is a segment register with no specific processor-defined purpose and is given purpose by the OS. In Linux, FS register is for thread local storage. We can also see the system call `access("/etc/ld.so.preload", R_OK)`, which tries to access `/etc/ld.so.preload` file. This file contains ELF shared objects, which are loaded before the execution of the program [7]. From the printout, we can see that this file is accessed before `libc.so` is loaded.

## 6 Measuring memory access behavior with background activity

In this section, we study the memory access behavior when there is a background activity. We use another process to allocate and access a large amount of memory to simulate background activity. We use `compete_for_memory` function for this purpose.

```

1 int compete_for_memory(void *unused) {
2     long mem_size = get_mem_size();
3     int page_sz = sysconf(_SC_PAGE_SIZE);
4     printf("Total memsize is %.3f GBs\n",
5           (double)mem_size / (1024 * 1024 * 1024));
6     fflush(stdout);
7     char *p = mmap(NULL, mem_size, PROT_READ | PROT_WRITE,
8                   MAP_NORESERVE | MAP_PRIVATE | MAP_ANONYMOUS, -1, 0);
9     if (p == MAP_FAILED) {
10         perror("Failed anon MMAP competition");
11         fflush(stdout);
12         exit(EXIT_FAILURE);
13     }
14
15     int i = 0;

```

---

<sup>9</sup>we trace our program as `strace ./memory_access 0 0 0 0 0 0`, where `memory_access` is the executable from `memory_access.c`

```

16     while (1) {
17         volatile char *a;
18         long r = simplerand() % (mem_size / page_sz);
19         char c;
20         if (i >= mem_size / page_sz) {
21             i = 0;
22         }
23         a = p + r * page_sz;
24         c += *a;
25         if (i % 8 == 0) {
26             *a = 1;
27         }
28         i++;
29     }
30     return 0;
31 }

```

`get_mem_size` is used to obtain the total amount of RAM of current system in bytes by reading from `/proc/meminfo` file. In the function above, `fflush(stdout)` is invoked after printing total memory. This is necessary because otherwise we cannot see the total memory message. However, flush to `stderr` is not necessary because each print call to `stderr` will be automatically followed by a flush. Since `compete_for_memory` function will try to allocate a memory region equals to the size of RAM, we need to setup swap file for the system because otherwise, the process will be killed due to oversubscribe of the memory. In the experiment, we setup a swap file of 2GB in the VM. `compete_for_memory` is called by a separate child process via `fork`. The child process will be killed when the main process exits. After forking the child process, the main process will sleep for 5 seconds to wait for child process to allocate memory. The statistics is shown in Figure 7. Note that the statistics is for anonymous mapping with random memory access. For other mapping and access behaviors, we submit the raw data along with the report. With memory pressure in the system, we expect the program takes more time to finish and encounters more page faults. From the statistics, we see that the program with memory pressure takes more time to finish and page faults are significantly higher than the no-compete case (shown both from `minflt` and `majflt`). However, read miss for both L1 cache and data TLB in compete case is slightly higher than the no-compete case. We also see there are non-zero values of `inblock` in the case of memory pressure. This is another indicator that more page faults happened. Other mappings (file-backed in shared and private mode) has similar statistics with slight higher page faults.

Next, we modify the page replacement policy of Linux kernel. Linux kernel uses LRU to reclaim pages. We modify `shrink_page_list` function in `mm/vmscan.c` file. The target code block looks like below:

```

1 switch (references) {
2 case PAGEREF_ACTIVATE:
3     goto activate_locked;
4 case PAGEREF_KEEP:
5     nr_ref_keep++;
6     goto keep_locked;
7 case PAGEREF_RECLAIM:
8 case PAGEREF_RECLAIM_CLEAN:
9     ; /* try to reclaim the page below */
10 }

```

We modify the `PAGEREF_ACTIVATE` case to behave the same as `PAGEREF_RECLAIM` case. Before the modification, `PAGEREF_ACTIVATE` represents the case where the kernel sees the page has been recently accessed. Thus, the page will not be swapped out of RAM. If we change this case to behave the same as `PAGEREF_RECLAIM`, there is no LRU in the page replacement process. We conduct our experiment for memory random access under memory pressure for all possible mappings. The statistics is shown in Figure 8. As one can see, since we disable the LRU replacement policy, the page faults increase in the modified kernel case. This is intuitive as the kernel might swap out pages that are frequently accessed. However, we should note that the page faults do not increase dramatically as in theory LRU can behave worse than the no replacement policy at all especially when the memory access pattern is random. This is because random memory access may create some access pattern that is in favor of no replacement policy at all case. For example, there is a clear page faults increase for file-backed mapping in private mode. However, for anonymous mapping, page faults actually

	[with compete_for_memory]	[no compete_for_memory]
uptime:	6.811 (std 0.164) s	4.790 (std 0.047) s
stime:	0.101 (std 0.130) s	0.000 (std 0.001) s
minflt:	19.800 (std 3.059)	10.200 (std 0.980)
majflt:	16.200 (std 2.482)	0.000 (std 0.000)
inblock:	1737.600 (std 734.998)	0.000 (std 0.000)
oublock:	0.000 (std 0.000)	0.000 (std 0.000)
maxrss:	674.400 (std 121.342)	0.000 (std 0.000)
Data L1 read access:	7521457650.400 (std 3348.175)	7521437569.600 (std 343.389)
Data L1 write access:	1077936131.000 (std 0.000)	1077936131.000 (std 0.000)
Data L1 read miss:	849878078.000 (std 3046793.999)	816180244.600 (std 1026125.094)
Data L1 read miss rate:	0.113 (std 0.000)	0.109 (std 0.000)
Data TLB read miss:	1660888.600 (std 24995.351)	1495066.600 (std 3024.350)
Data TLB write miss:	779031.000 (std 11851.924)	715924.000 (std 1492.899)

Figure 7: The statistics of anonymous mapping for random memory access both with and without memory pressure



[Anonymous mapping]		
	[modified kernel]	[no modified kernel]
utime:	4.824 (std 0.007) s	4.773 (std 0.009) s
stime:	0.001 (std 0.001) s	0.001 (std 0.001) s
minflt:	24.200 (std 1.600)	22.600 (std 1.020)
majflt:	3.400 (std 1.744)	4.200 (std 0.400)
inblock:	500.800 (std 250.879)	624.000 (std 16.000)
oublock:	0.000 (std 0.000)	0.000 (std 0.000)
maxrss:	0.000 (std 0.000)	0.000 (std 0.000)
Data L1 read access:	7521438898.800 (std 127.804)	7521437634.000 (std 211.443)
Data L1 write access:	1077936131.000 (std 0.000)	1077936131.000 (std 0.000)
Data L1 read miss:	825365597.800 (std 1754742.247)	819683681.600 (std 381926.508)
Data L1 read miss rate:	0.110 (std 0.000)	0.109 (std 0.000)
Data TLB read miss:	1509795.600 (std 4720.234)	1482084.800 (std 822.343)
Data TLB write miss:	723814.600 (std 2037.219)	712125.200 (std 1102.349)
[File-based mapping (private)]		
	[modified kernel]	[no modified kernel]
utime:	4.844 (std 0.005) s	4.784 (std 0.008) s
stime:	0.000 (std 0.001) s	0.001 (std 0.001) s
minflt:	18.200 (std 1.470)	22.000 (std 1.549)
majflt:	10.600 (std 0.800)	4.200 (std 1.600)
inblock:	2212.800 (std 250.930)	606.400 (std 291.200)
oublock:	0.000 (std 0.000)	0.000 (std 0.000)
maxrss:	149.600 (std 299.200)	170.400 (std 340.800)
Data L1 read access:	7521438566.000 (std 368.873)	7521437760.800 (std 346.232)
Data L1 write access:	1077936131.000 (std 0.000)	1077936131.000 (std 0.000)
Data L1 read miss:	828836805.000 (std 987778.815)	822928672.600 (std 1469473.160)
Data L1 read miss rate:	0.110 (std 0.000)	0.109 (std 0.000)
Data TLB read miss:	1521589.600 (std 2499.976)	1493323.200 (std 2975.783)
Data TLB write miss:	729151.200 (std 1631.806)	718707.400 (std 1034.893)
[File-based mapping (shared)]		
	[modified kernel]	[no modified kernel]
utime:	4.868 (std 0.032) s	4.792 (std 0.019) s
stime:	0.015 (std 0.015) s	0.014 (std 0.013) s
minflt:	4125.800 (std 1.166)	4126.000 (std 1.265)
majflt:	3.800 (std 0.400)	4.000 (std 0.000)
inblock:	593.600 (std 44.800)	616.000 (std 0.000)
oublock:	32832.000 (std 0.000)	32832.000 (std 0.000)
maxrss:	2843.200 (std 167.033)	2466.400 (std 726.388)
Data L1 read access:	7521441059.200 (std 1197.569)	7521440534.000 (std 760.447)
Data L1 write access:	1077936131.000 (std 0.000)	1077936131.000 (std 0.000)
Data L1 read miss:	829703258.400 (std 2059767.610)	823238818.400 (std 691773.770)
Data L1 read miss rate:	0.110 (std 0.000)	0.109 (std 0.000)
Data TLB read miss:	1545167.800 (std 7327.071)	1494720.400 (std 1412.020)
Data TLB write miss:	744986.600 (std 3014.310)	718793.000 (std 1569.477)

Figure 8: The statistics of memory random access behavior under different mappings before and after the modification of the kernel

decrease after the modification. In addition, since the `do_mem_access` frequently accesses around 16MB of memory, while the competing process uniformly accesses its 2GB memory, the chance that a hot page being swapped out is relative small.

## References

- [1] “script(1) - linux man page.” <http://man7.org/linux/man-pages/man1/script.1.html>.
- [2] “strace(1) - linux man page.” <https://linux.die.net/man/1/strace>.
- [3] “lsof(8) - linux man page.” <http://man7.org/linux/man-pages/man8/lsof.8.html>.
- [4] “Dynamic host configuration protocol.” [https://en.wikipedia.org/wiki/Dynamic\\_Host\\_Configuration\\_Protocol](https://en.wikipedia.org/wiki/Dynamic_Host_Configuration_Protocol), 2018.
- [5] “Hardware performance counter.” [https://en.wikipedia.org/wiki/Hardware\\_performance\\_counter](https://en.wikipedia.org/wiki/Hardware_performance_counter), 2018.
- [6] “arch\_prctl(2) - linux man page.” [http://man7.org/linux/man-pages/man2/arch\\_prctl.2.html](http://man7.org/linux/man-pages/man2/arch_prctl.2.html).
- [7] “ld.so(8) - linux man page.” <http://man7.org/linux/man-pages/man8/ld.so.8.html>.