

# CS380L: Advanced Operating Systems Lab #1

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

EID:zh4378 Spring 2019

## 1 Environment

Unless otherwise noted, we use a Linux server for all the experiments. The server has 4 Intel(R) Xeon(R) CPU E3-1220 v5 @ 3.00GHz processors and 16GB of memory, and runs Ubuntu 16.04.2 LTS (kernel version 4.11.0). The CPU has 32KB of L1 data cache per core (8-way set associative) (found through `getconf -a | grep CACHE`). In addition, it has two-level TLBs. The first level (data TLB) has 64 entries (4-way set associative), and the second level has 1536 entries for both instructions and data (6-way set associative) (found through `cpuid | grep -i tlb`).

## 2 Memory map

`/proc/[pid]/maps` file contains process `[pid]`'s mapped memory regions and their access permissions [1]. We use the following code to read content of `/proc/self/maps` file <sup>2</sup>:

```
1 sprintf(filepath, "/proc/%u/maps", (unsigned) getpid());
2 FILE *f = fopen(filepath, "r");
3
4 printf("%-32s %-8s %-10s %-8s %-10s %s\n", "address", "perms", "offset", "dev", "
   inode", "pathname");
5 while (fgets(line, sizeof(line), f) != NULL) {
6     sscanf(line, "%s%s%s%s%s", address, perms, offset, dev, inode, pathname);
7     printf("%-32s %-8s %-10s %-8s %-10s %s\n", address, perms, offset, dev, inode,
   pathname);
8 }
9 fclose(f);
```

In the file, each line corresponds to a mapped memory region. There are six columns of each line, which represent six properties of the mapped memory region: `address`, `perms`, `offset`, `dev`, `inode`, and `pathname`. The result of running `memory_map.c` is below:

The `address` field gives range of virtual memory address of the mapped memory region. Access permission of each memory region is indicated by `perms` field. There are four bits in the field: `rw` `x`

---

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

<sup>2</sup>see `memory_map.c` for complete code

address	perms	offset	dev	inode	pathname
00400000-00401000	r-xp	00000000	fd:01	12374202	/home/zeyuanhu/380L-Spring19/lab1/src/a.out
00600000-00601000	r--p	00000000	fd:01	12374202	/home/zeyuanhu/380L-Spring19/lab1/src/a.out
00601000-00602000	rw-p	00001000	fd:01	12374202	/home/zeyuanhu/380L-Spring19/lab1/src/a.out
022c1000-022e2000	rw-p	00000000	00:00	0	[heap]
7fea45315000-7fea454d5000	r-xp	00000000	fd:01	24903836	/lib/x86_64-linux-gnu/libc-2.23.so
7fea454d5000-7fea456d5000	---p	001c0000	fd:01	24903836	/lib/x86_64-linux-gnu/libc-2.23.so
7fea456d5000-7fea456d9000	r--p	001c0000	fd:01	24903836	/lib/x86_64-linux-gnu/libc-2.23.so
7fea456d9000-7fea456db000	rw-p	001c4000	fd:01	24903836	/lib/x86_64-linux-gnu/libc-2.23.so
7fea456db000-7fea456df000	rw-p	00000000	00:00	0	/lib/x86_64-linux-gnu/libc-2.23.so
7fea456df000-7fea45705000	r-xp	00000000	fd:01	24903834	/lib/x86_64-linux-gnu/ld-2.23.so
7fea458e0000-7fea458e3000	rw-p	00000000	00:00	0	/lib/x86_64-linux-gnu/ld-2.23.so
7fea45904000-7fea45905000	r--p	00025000	fd:01	24903834	/lib/x86_64-linux-gnu/ld-2.23.so
7fea45905000-7fea45906000	rw-p	00026000	fd:01	24903834	/lib/x86_64-linux-gnu/ld-2.23.so
7fea45906000-7fea45907000	rw-p	00000000	00:00	0	/lib/x86_64-linux-gnu/ld-2.23.so
7ffe67d7e000-7ffe67d9f000	rw-p	00000000	00:00	0	[stack]
7ffe67da3000-7ffe67da5000	r--p	00000000	00:00	0	[vvar]
7ffe67da5000-7ffe67da7000	r-xp	00000000	00:00	0	[vdso]
fffffffff600000-fffffffff601000	r-xp	00000000	00:00	0	[vsyscall]

Figure 1: Output of memory\_map.c

represents read, write, and executable respectively; the last bit (p or s) represents whether the region is private or shared. **offset** field represents the offset in the mapped file. **dev** field indicates the device (represented with format of **major:minor**) that the mapped file resides . There are two kinds of value in this column for our case: **fd:01** and **00:00**. The former one is the device id (in hex) of / (checked with `mountpoint -d /`) and the latter one represents no device associated with the file. **inode** field represents the inode number of the file on the device. 0 means no file is associated with the mapped memory region. **pathname** field gives the absolute path to the file associated with the mapped memory region. It can be some special values like `[heap]`, `[stack]`, `[vdso]`, etc.

To locate the start of the text section of the executable, we invoke `objdump -h` on the binary and get `0000000000400600`. Output of `/proc/self/maps` shows that the start address of `libc` is `7fea45315000`. The reason for these two addresses are different is `libc` is dynamic loaded library, which is loaded during the runtime of executable, which is not compiled and linked as part of executable. The code segment contains the executable instruction, not the dynamic loaded library.

One interesting thing happens between runs of the executable: the content of `/proc/self/maps` is different. Addresses of all mapped memory regions are different except for the regions mapped to the executable and `[vsyscall]`. The root cause behind this phenomenon is Address Space Layout Randomization (ASLR) [2] for programs in user space. This feature is enabled by default and can be seen via the content of `/proc/sys/kernel/randomize_va_space` file. In our case, the value is 2, which means the positions of stack itself, virtual dynamic shared object (VDSO) page, shared memory regions, and data segments are randomized [3].

### 3 Getrusage

To get resource usage of the current process, we use `getrusage` [4]. The result is stored in `rusage` struct. Not all fields of the struct are completed: unmaintained fields are set to zero by the kernel. Those fields exist for compatibility with other systems purpose. The following code instantiates `rusage` struct and print all the maintained fields <sup>3</sup>:

```
1 struct rusage usage;
2 if (getrusage(RUSAGE_SELF, &usage) != 0) {
3     perror("getrusage");
4     return 0;
5 }
6
7 // user CPU time used
8 printf("utime = %ld.%06ld s\n", usage.ru_utime.tv_sec,
9      usage.ru_utime.tv_usec);
10 // system CPU time used
11 printf("stime = %ld.%06ld s\n", usage.ru_stime.tv_sec,
12      usage.ru_stime.tv_usec);
13 // maximum resident set size
14 printf("maxrss = %ld KB\n", usage.ru_maxrss);
15 // page reclaims (soft page faults)
16 printf("minflt = %ld\n", usage.ru_minflt);
17 // page faults (hard page faults)
18 printf("majflt = %ld\n", usage.ru_majflt);
19 // block input operations
20 printf("inblock = %ld\n", usage.ru_inblock);
21 // block output operations
22 printf("oublock = %ld\n", usage.ru_oublock);
23 // voluntary context switches
24 printf("nvcsw = %ld\n", usage.ru_nvcsw);
25 // involuntary context switches
26 printf("nivcsw = %ld\n", usage.ru_nivcsw);
```

man page of `getrusage` explains the meaning of each field [4] in details. `utime` and `stime` are about CPU time usage; `minflt` and `majflt` are related to page faults; `maxrss` represents the maximum size of working set; `inblock` and `oublock` are about file system I/O.

---

<sup>3</sup>complete code can be seen in `getrusage.c`

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

## 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 [5] 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 [6]). `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`).

According to man page [5], 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),

```

0000000000400d36 <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

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 values. 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));
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) {

```

<sup>4</sup>complete code can be seen in `perf_event_open_usage.c`

```

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 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++) {
```

---

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



```

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()` 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()`).

	[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

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 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>.

<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

<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;`

	[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.

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 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` [7]. 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`

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

("/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 [8]. 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 %3.2f 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;
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.

	[with compete_for_memory]	[no compete_for_memory]
utime:	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

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 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] “proc(5) - linux man page.” <http://man7.org/linux/man-pages/man5/proc.5.html>.
- [2] “Address space layout randomization (aslr).” [https://en.wikipedia.org/wiki/Address\\_space\\_layout\\_randomization](https://en.wikipedia.org/wiki/Address_space_layout_randomization), 2018.



[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

- [3] “Linux and aslr: kernel/randomize\_va\_space.” [https://linux-audit.com/linux-aslr-and-kernelrandomize\\_va\\_space-setting](https://linux-audit.com/linux-aslr-and-kernelrandomize_va_space-setting), 2016.
- [4] “getrusage(2) - linux man page.” <http://man7.org/linux/man-pages/man2/getrusage.2.html>.
- [5] “perf\_event\_open(2) - linux man page.” [http://man7.org/linux/man-pages/man2/perf\\_event\\_open.2.html](http://man7.org/linux/man-pages/man2/perf_event_open.2.html).
- [6] “Hardware performance counter.” [https://en.wikipedia.org/wiki/Hardware\\_performance\\_counter](https://en.wikipedia.org/wiki/Hardware_performance_counter), 2018.
- [7] “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).
- [8] “ld.so(8) - linux man page.” <http://man7.org/linux/man-pages/man8/ld.so.8.html>.