

گام یک:

سوال ۱ و ۳: مود های گوناگون برای cpu حالت های اجرای اند که در هر یک از آنها محدودیت مشخصی روی دستوراتی که cpu و برخی سخت افزار های خاص می تواند اجرا کند قرار داده می شود. این نوع طراحی امکان ارائه خدمات اضافه به یک کاربر و یا در یک مود جدید در مقایسه با بقیه را امکان پذیر می سازد. در تئوری فقط دستورات اصلی کرنل در حالت کرنل اجرا می شود و بقیه برنامه ها با استفاده از سیستم کال ها با سیستم عامل در ارتباط خواهند بود و دستورات سطح سیستم را اجرا می کنند. در عمل چون سیستم کال زمان می برد و سربار دارد پس برخی برنامه ها که حساسیت به زمان دارند (مثلا device driver ها) با دسترسی کرنل مود اجرا می شوند.

کرنل مود حالتی است در که در آن کد به تمام زیر ساخت های سخت افزاری دسترسی کامل دارد، می تواند به هر آدرس خانه ای حافظه دسترسی داشته باشد و هر دستوری را که می خواهد در CPU اجرا کند. در حالت یوزر کد امکان دسترسی مستقیم به سخت افزار و حافظه را ندارد. همچنین امکان دسترسی به دستگاه های ورودی و خروجی را نداریم و به طور کلی در حالت یوزر دستوراتی که می تواند حالت عمومی سیستم را تغییر دهد و روی سایر پروسس ها تاثیر بگذارند قابل اجرا نیستند. بخش عمده ای از برنامه هایی که ما به طور روزانه اجرا می کنیم در حالت یوزر هستند. اگر می خواهیم دسترسی های کرنل را به تمام یوزرها بدهیم امنیت سیستم قابل تامین نبود.

سوال ۲: دسته بندی فراخوانی های سیستمی با توجه به کارکردهایشان:

Process control

end, abort
load, execute
create process, terminate process
get process attributes, set process attributes
wait for time
wait event, signal event
allocate and free memory

File management

create file, delete file
open, close file
read, write, reposition
get and set file attributes

Device management

request device, release device
read, write, reposition
get device attributes, set device attributes
logically attach or detach devices

Information maintenance

get time or date, set time or date
get system data, set system data
get and set process, file, or device attributes

Communications

create, delete communication connection
send, receive messages
transfer status information
attach and detach remote devices

مثال هایی در ویندوز و یونیکس:

Process control: unix: fork(), exit(), wait()
Windows: CreateProcess(), ExitProcess(), WaitForSingleObject()
File management: unix: open(), read(), write(), close()
Windows: CreateFile(), ReadFile(), WriteFile(), CloseHandle()
Device management: unix: ioctl(), read(), write()
Windows: SetConsoleMode(), ReadConsole(), WriteConsole()
Information maintenance: unix: getpid(), alarm(), sleep()
Windows: GetCurrentProcessID(), SetTimer(), Sleep()
Communication: unix: pipe(), shmget(), mmap()
Windows: CreatePipe(), CreateFileMapping(), MapViewOfFile()
Protection: unix: chmod(), umask(), chown()
Windows: SetFileSecurity(), InitializeSecurityDescriptor(), SetSecurityGroup()

فراخوانی سیستمی دو بخش دارد:

برای فضای کاربر امکان پذیر نیست که کد kernel را به طور مستقیم اجرا کند. برنامه های فضای کاربر باید به kernel یک سیگنال ارسال کنند که میخواهند یک system call را اجرا کنند و سپس به حالت کرنل سوییچ کنند. مکانیزم فرستادن سیگنال به کرنل یک intrupt نرم افزاری است. یک exception ایجاد میشود؛ سیستم به kernel mode تغییر میکند و exception handler را اجرا میکند. در اینجا exception handler همان system call handler است. این intrupt نرم افزاری در x86، داخل 128 intrupt number است. که از طریق دستور 0x80\$ انجام میشود. این مساله باعث سوییچ به کرنل و اجرای exception vector 128 که همان system call handler است میشود. وارد شدن به فضای کرنل به تنهایی کافی نیست زیرا فراخوانی های سیستمی چندگانه ای وجود دارند که همه به شیوه ی مشابه وارد کرنل میشوند. بنابراین برای تشخیص نوع فراخوان سیستمی باید system call number به کرنل فرستاده شود. این مقدار در رجیستر eax ذخیره میشود. System call handler مقدار eax را میخواند. تابع system-call() اعتبار system call number را میسنجد و در صورت نامعتبر بودن تابع ENO-sys را باز میگرداند. سپس آرگومان های فراخوانی سیستمی از طریق رجیسترهای سخت افزاری دیگر منتقل شده و با اجرای یک دستور intrupt صورت میگیرد. بعد از اجرای intrupt، شماره ی فراخوانی سیستمی به عنوان ایندکس جدولی از اشاره گر ها استفاده میشود تا آدرس شروع کد راه انداز فراخوانی سیستمی به دست آید. سپس process به این آدرس پریده و امتیاز های process از مد کاربر به مد هسته منتقل می شوند. فهرست شماره های فراخوانی های سیستمی نسخه های اخیر هسته ی لینوکس در مسیر include/asm-x86/unistd.h/ موجود است. (برای نمونه، NR_close، که متناظر با فراخوانی سیستمی close() است برای بستن یک فایل استفاده میشود و با مقدار ۶ تعریف میشود) معمولاً فهرست اشاره گر های راه انداز فراخوانی های سیستمی در فایل arch/x86/kernel/entry-64.s زیر سر خط ENTRY(sys_call_table) میشود.

entry_64.S:

- * Syscall return path ending with SYSRET (fast path)
- * We must check ti flags with interrupts (or at least preemption)
- * off because we must *never* return to userspace without
- * processing exit work that is enqueued if we're preempted here.
- * In particular, returning to userspace with any of the one-shot
- * flags (TIF_NOTIFY_RESUME, TIF_USER_RETURN_NOTIFY, etc)
- * Has incomplete stack frame and undefined top of stack
- * Reload arg registers from stack in case ptrace changed them.
- * We don't reload %rax because syscall_trace_entry_phase2() returned
- * the value it wants us to use in the table look
- * edi: mask to check *
- * edx: work, edi: workmask
- handle signals and tracing — both require a full stack frame
- * Check for syscall exit trace */

منابع:

http://www.cs.iit.edu/~cs561/cs450/system_calls/style/9.html <http://scanfreetree.com/operating-system/Types-of-System-Calls>

آیا روش های دیگری غیر از فراخوانی سیستمی برای دسترسی به اطلاعات و سرویس های هسته وجود دارد؟

در لینوکس system call ها ابزار هایی هستند که interface بین فضای کاربر و هسته هستند. system call ها تنها ورودی مجاز به هسته هستند به غیر از exception ها و trap ها. در واقع interface های دیگر از قبیل device files یا proc/ نهاییات از طریق system call ها قابل دسترسی هستند. یکی از روش های دسترسی به سرویس های کرنل استفاده از API ها هستند. در واقع با صدا زدن یک API یک سیستم کال فراخوانی میشود. علت استفاده از API درگیر نشدن با جزئیات و portability بودن برنامه است. یک روش دیگر message passing است. در این روش process ایجاد شده در فضای کاربر یک message ایجاد میکند و سرویس مورد در خواستش را توصیف میکند. بعد از آن از یک تابع مطمئن برای فرستادن message به یک process مورد اعتماد سیستم عامل استفاده میکند. این تابع مشابه یک trap عمل میکند و message را بررسی کرده و به کرنل مود سوییچ میکند. به همین روش نتیجه ی سرویس در خواستی را به فضای کاربر برمیگرداند.

منبع: کتاب love

سوال ۴: یکی از روش ها انتقال پارامترها با استفاده از رجیستر است که در روش تعداد پارامترهایی که می توانیم منتقل کنیم محدود است و اینکه دقیقاً از کدام رجیستر ها برای انجام این کار می توانیم استفاده کنیم به ماشین ما بستگی دارد مثلاً:

On [x86-64 UNIX systems](#), including Linux and default NetRun, the first six parameters go into rdi, rsi, rdx, rcx, r8, and r9.

On [Windows 64](#), the first four parameters go into rcx, rdx, r8, and r9

On 32-bit x86 systems, no parameters go in registers, they're all on "the stack" (wait a month to hear about this!).

پس زمانی که می خواهیم تعدادی بیشتر از ۶ پارامتر (برای یونیکس مثلاً) پاس دهیم به سراغ استک می رویم.

<https://stackoverflow.com/questions/19858980/how-does-assembly-do-parameter-passing-by-value-reference-pointer-for-differe>

سوال ۵: در رجیستر ها ذخیره می شوند. در x86-32 ، رجیستر های edi , ebx , ecx , edx , esi به ترتیب حاوی 5 آرگومان اول می باشند. در صورت بیشتر بودن تعداد آرگومان ها، یک رجیستر برای نگهداری اشاره گر به فضای کاربر نگهداری می شود که در آن فضا، همه ی پارامتر ها ذخیره شده اند. مقدار بازگشتی هم با رجیستر به فضای کاربر داده می شود. که در x86، در رجیستر eax نوشته می شود.

منبع: بخش parameter passing کتاب love صفحه 74

سوال 6: الف) کرنل دو متد برای بررسی های لازم و کپی کردن از فضای کاربر و کپی کردن به آن تهیه می کند. برای نوشتن در فضای کاربر، متد copy_to_user () استفاده می شود که 3 پارامتر می گیرد. پارامتر اول آدرس حافظه مقصد در فضای آدرس process می باشد و پارامتر دوم، اشاره گر به مبدایی است در فضای کرنل که می خواهیم از آن کپی کنیم، و پارامتر سوم هم سایز داده به بایت می باشد. برای خواندن از فضای کاربر، از متد copy_from_user () استفاده می شود که شبیه copy_to_user () می باشد. این تابع از پارامتر دوم به اندازه ی پارامتر سوم بایت داده میخواند و در پارامتر اول می ریزد. هر دو این تابع ها تعداد بایت هایی را که نتوانسته اند بخوانند را برمی گردانند و در صورت موفقیت صفر برمی گردانند.

ب) حافظه به دو بخش مجزا تقسیم می شود. 1) فضای کاربر که قسمتی که پروسه های کاربر در آن اجرا می شوند. نقش کرنل مدیریت برنامه های در حال اجرا در این فضا می باشد که با هم تداخل نداشته باشند. 2) فضای کرنل که فضایی است که کد کرنل در آن ذخیره و اجرا می شود. پروسه های در حال اجرا در فضای کاربر فقط به بخش محدودی از حافظه دسترسی دارند، در حالی که کرنل به تمام حافظه دسترسی دارد. علاوه بر این، پروسه های فضای کاربر به فضای کرنل دسترسی ندارند و فقط می توانند توسط رابطی که کرنل در اختیار آنها قرار می دهد (system calls) به بخش کوچکی از کرنل دسترسی پیدا کنند.

برنامه های فضای کاربر و کرنل، فضای آدرس های متفاوتی دارند به همین دلیل کپی کردن از یکی به دیگری نیازمند تغییر فضای آدرس میباشد.

<https://unix.stackexchange.com/questions/87625/what-is-difference-between-user-space-and-kernel-space>

SYSCALL:

یک تابع کتابخانه ای است که system call را فراخوانی میکند که interface زبان اسمبلی آن عدد و آرگومان خاصی دارد. استفاده از syscall مفید است به عنوان مثال برای system call ای که wrapper function در کتابخانه ی c ندارد. (syscall) رجیسترهای cpu را قبل از ایجاد system call ذخیره میکند. رجیستر ها را بعد از بازگشت از system call بازیابی میکند. و هر error code ای که از system call باز می گردد را در errno ذخیره میکند. (در صورت رخ دادن error)

منبع: <http://man7.org/linux/man-pages/man2/syscall.2.html>

گام دو:

اضافه کردن یک سیستم کال برای هر سرویس جدیدی که نیاز داریم منطقی نیست زیرا برای هر سیستم کالی که اضافه کنیم لازم است که کرنل را مجدداً make کنیم که این کار زمان بسیار زیادی طول می‌کشد از طرفی نیز تعدد سیستم‌کال‌های غیر استاندارد باعث می‌شود که قابلیت انتقال و portability کاهش پیدا کند چون آن سیستم کال‌ها روی سیستم‌های عادی دیگر نخواهد بود پس برنامه ما اجرا نخواهد شد.

برای اضافه کردن هر سیستم کال یک فولدر در پوشه‌ای که کد کرنل ما در آن قرار دارد ایجاد کردیم. سپس در داخل آن پوشه یک فایل C ایجاد کردیم و کد مورد نظرمان را در آن نوشتیم. سپس یک make file در همین فولدر ایجاد کرده و عبارت `obj-y := نام فایل حاوی کد.o` را در آن نوشتیم. پس از آن به فولدر حاوی source کرنل بازگشتیم و فایل `make file` را در آن باز کردیم. در خط ۸۴۲ آن عبارت `core-y += kernel/ mm/ fs/ ipc/ security/ crypto/ block/ name of folder` را در آن تغییر می‌دهیم. سپس به بخش `cd arch/x86/syscalls` رفته و فایل `syscall_64.tbl` را باز می‌کنیم و عبارتی به شکل زیر را به انتهای فایل اضافه می‌کنیم: `number of last system call+1 64 name of file sys_name of file`. سپس به بخش `cd include/linux` رفته و فایل `syscalls.h` را باز می‌کنیم و به آن `asm linkage long name of function` اضافه می‌کنیم. سپس کرنل را کامپایل می‌کنیم پس از اتمام آن یک کد تست برای فراخوانی سیستمی جدید اضافه شده می‌نویسیم و آن را کامپایل می‌کنیم. در صورت اضافه شدن صحیح سیستم کال با وارد کردن دستور `dmesg` می‌توانیم اجرای آن را ببینیم. کدهای مربوط به فراخوانی‌های سیستمی مورد سوال و تست آن‌ها به همراه فایل‌هایی که با تغییرات همراه بوده‌اند (`Make file, syscalls.h, syscall_64.tbl`) ضمیمه شده است.

گام سوم:

سوال ۱: سیستم کالی با نام ptrace که کوتاه شده‌ی process trace است در لینوکس وجود دارد (<https://linux.die.net/man/2/ptrace>). این سیستم کال امکانی را فراهم می‌کند که یک پردازنده بتواند پردازنده دیگری را ردگیری کند و آن را بررسی کرده و حتی می‌تواند حافظه‌ی آن را هم تغییر دهد. این سیستم کال اغلب برای دیباگ کردن در حالتی که breakpoint قرار داده شده و به طور کلی ردگیری سیستم کال‌ها استفاده می‌شود. برای استفاده از ptrace باید اول به آن پردازنده ای که می‌خواهیم ردگیری کنیم متصل شویم و این کار به صورت ۱ thread ۱ thread انجام می‌شود. تعریف کلی این سیستم کال به شکل زیر است:

```
#include <sys/ptrace.h>
```

```
long ptrace(enum __ptrace_request request, pid_t pid,  
void *addr, void *data);
```

که در آن `pid` برابر با شناسه‌ی آن thread ای است که می‌خواهیم ردگیری کنیم. مقدار بازگشتی این تابع `status` را باز می‌گرداند که مشخص می‌کند پردازنده یا ریسمان در حال ردگیری در چه وضعیتی است. `request` مشخص می‌کند که چه کاری می‌خواهیم انجام دهیم مثلاً اگر مقدار آن `PTRACE_TRACEME` باشد مشخص می‌کند که پردازنده باید توسط پدرش ردگیری شود در این نوع درخواست سایر پارامترهای ورودی نادیده گرفته می‌شوند. ولی مثلاً اگر مقدار این درخواست برابر با `PTRACE_POKETEXT, PTRACE_POKEDATA` باشد داده‌ی `data` را در آدرس `addr` از حافظه‌ی پردازنده در حال ردگیری می‌نویسد و یا دستور `PTRACE_GETREGSET` مقدار رجیسترهای پردازنده در حال ردگیری را تغییر می‌دهد. که `addr` به آن رجیستر اشاره کرده و `data` اندازه و آدرس بافری است که داده‌ای که می‌خواهیم در رجیستر بنویسیم در آن قرار دارد. می‌تواند مقادیر بسیار متنوع دیگری هم داشته باشد.

دستور ptrace باعث می‌شود که پردازنده‌ای که در حال ردگیری آن است هر بار که با یک سیستم کال به کرنل داخل شده و یا از آن خارج می‌شود متوقف شود تا بتوان با روی آن تغییرات لازم را داد و یا حالت کلی آن را بررسی کرد و بعد ادامه می‌دهد و این روند تکرار می‌شود.

`strace` برای ردگیری کردن سیستم کال‌ها و سیگنال‌ها استفاده می‌شود، در حالت ساده این دستور سیستم کال‌های استفاده شده توسط یک پردازنده، نام‌های آن‌ها و سیگنال‌های بازگشتی آن‌ها به پردازنده را ردگیری می‌کند. می‌توان از ptrace برای پیاده‌سازی آن استفاده کرد چون ما در `strace` می‌خواهیم که پارامترهای داده شده به سیستم کال و سیستم کال صدا شده را ببینیم همان طور که بالاتر گفتیم ptrace پس از صدا شدن سیستم کال می‌ایستد و مثلاً در مدل x86 پارامترها با استفاده از رجیسترها پاس داده می‌شوند که می‌توانیم با استفاده از درخواست `PTRACE_GETREGS` از ptrace مقدار رجیسترها را کپی کنیم. و برای فهمیدن اینکه کدام سیستم کال هم صدا شده، هر سیستم کال یک عدد خاص خود را دارد که توسط رجیسترها فرستاده می‌شود پس با خواندن مقدار آن‌ها نوع سیستم کال را هم متوجه می‌شویم و همان طور که گفتیم پس از پایان سیستم و بازگشت از کرنل هم ptrace متوجه شده و پردازنده را متوقف می‌کند پس باز هم `strace` می‌تواند مقدار بازگشتی که همان سیگنال فرستاده شده به پردازنده است را ردگیری کند.

سوال ۲: دلیل استفاده از کتابخانه‌ها به جای دسترسی مستقیم این است که اگر هر برنامه نویس بخواهد به طور مستقیم از خود سیستم کال‌ها استفاده کند پس تمام برنامه نویسان باید شهود خوبی نسبت به سیستم و زیرساخت‌های آن و نحوه‌ی پیاده‌سازی سیستم کال‌ها و سیستم کال‌های موجود داشته باشند که این خود کار دشواری است. از طرفی نیز چون سیستم کال‌ها الزاماً بین تمام سیستم‌ها یکی نیستند و گاهی وابسته به ماشین و یا سیستم عامل اند اگر بخواهیم در برنامه‌های خود به طور مستقیم

استفاده کنیم برنامه **portability** خود را از دست خواهد داد. پس ترجیح می‌دهیم به جای اینکه خود را نحوه‌ی پیاده‌سازی سطح پایین دستورات درگیر کنیم و **portability** برنامه خود را کم کنیم از کتابخانه‌های استاندارد استفاده کنیم تا برنامه روی تمام سیستم‌ها قابل اجرا باشد و نوشتن آن هم ساده‌تر باشد

سوال ۳:

```
void* malloc (size_t size);  
void *sbrk(intptr_t increment);  
void *mmap(void *addr, size_t length, int prot, int flags, int fd, off_t offset)
```

این تابع به اندازه‌ی **size** بایت در حافظه فضا می‌گیرد و اشاره‌گری به آن را برمی‌گرداند. در پیاده‌سازی این تابع از ۲ سیستم کال استفاده شده است. **Mmap** وظیفه‌ی **map or unmmap** کردن در حافظه را دارد. سیستم کال دوم **sbrk** است که اندازه‌ی بخش داده‌ی برنامه‌ی ما را تغییر می‌دهد. به طور دقیق‌تر به بررسی **Mmap** می‌پردازیم. (<https://linux.die.net/man/2/mmap>) این دستور یک تناظر جدید در حافظه‌ی مجازی پردازنده ایجاد می‌کند که شروع این تناظر از **addr** بوده و آدرس تناظر جدید را به عنوان مقدار برگشتی به ما باز می‌گرداند. **prot** نوع محافظت از حافظه جدید را مشخص می‌کند که آیا می‌تواند خوانده شود و یا در آن نوشته هم بشود و یا ... **flags** مشخص می‌کند که آیا به روز رسانی‌ها در این تناظر به اطلاع سایر پردازنده‌هایی که از آن استفاده می‌کنند هم برسد و یا خیر. **Fd** برابر با **file descriptor** است و تناظر ما به اندازه‌ی **length** و با شروع از **offset** انجام خواهد شد. شکل زیر برای کمک به فهم بهتر منظور از تناظر در حافظه مجازی پردازنده است.

به سراغ **sbrk** می‌رویم. (<https://linux.die.net/man/2/sbrk>) این دستور **program break** که به پایان بخش داده‌ی برنامه اشاره می‌کند را تغییر می‌دهد. این سیستم کال اندازه‌ی اشاره‌کننده به انتهای بخش داده را تا زیاد می‌کند و در صورت موفقیت مقدار قبلی **program break** که در واقع شروع حافظه‌ی گرفته شده‌ی جدید است را باز می‌گرداند و در صورت عدم موفقیت مقدار ۱- برمی‌گرداند. استفاده‌ی مستقیم از این دستور به جای **malloc** پیشنهاد نمی‌شود چون پارامترهای آن برای سیستم‌عامل‌های مختلف متفاوت اند و این قابلیت انتقال برنامه را پایین می‌آورد.