

# سیستم‌های عامل پیشرفته – تمرین یک

استاد: دکتر حسین اسدی

محمد مهدی قره‌گوزلو – 401206836

**پرسش یکم:** با توجه به اینکه هر ترد 1 میلیون بار 1024 کیلوبایت حافظه درخواست می‌دهد یعنی 1 ترابایت فقط توسط یک ترد درخواست داده می‌شود، با توجه به تعداد تردها در مجموع 256 ترابایت حافظه می‌شود (کلا 256 گیگابایت حافظه اصلی داریم).

فرض کنیم کلا 3 نوع shutdown داشته باشیم: local با دانه بندی ریز، ipi یا همان remote با دانه بندی ریز و Flush که تنها زمان تعویض متن رخ دهد. در این مورد بیشترین shutdown ها احتمالا از نوع remote ipi خواهند بود.

اگر 160 ترد روی هسته‌های مختلف در حال اجرا باشند برای هر صفحه که تخصیص می‌دهند یا آزاد می‌کنند باید یک ipi صادر شده و به هسته‌های دیگر این تغییر اعلام شود. با توجه به اندازه صفحه سیستم مورد نظر این تعداد ipi ها می‌تواند تغییر کند. مثلا چیزی در حدود 250 هزار ipi به ازای هر ترد هنگام تخصیص حافظه و همین مقدار برای آزادسازی (1TB / 4KB). واقعا عدد مشخصی نمی‌توان تعیین کرد، حتی در عمل نیز این اعداد نوسان دارند و صرفا در حد تخمین می‌توان نظر داد.

حالا اگر در این میان یکی از این تردها با تعویض متن از هسته خارج شود و اگر پردازنده مورد نظر ما شماره پردازنده‌ها را در TLB در کنار رکوردها نگه داری نکرده باشد یک flush کامل را نیز مجبوریم پردازیم.

بهتر است برای انجام عملیات‌هایی مانند تخصیص حافظه با تکرار بالا که ساختار جدول صفحه را با نوسان زیاد تغییر می‌دهند وقفه را غیر فعال کنیم به 2 دلیل: اول اینکه چون صرفا در حال تخصیص و آزاد کردن هستیم و نوشتن یا خواندن صورت نمی‌گیرد نیازی به shutdown فوری نداریم چون خطری از جهت خواندن یا نوشتن غیرمجاز یا نا معتبر وجود ندارد. توجه کنید فقط در صورتی کل عملیات‌ها تخصیص و آزاد کردن باشند. و دوم اینکه مشکل flush کردن نیز اگر وقفه غیرفعال شود و تعویض متن صورت نگیرد برطرف می‌شود. برای حالت دوم البته شاید بهتر باشد شماره pid ها را در TLB نگه داریم.

**توضیحاتی مختصر درباره دلایل Shutdown:** آنچه در ادامه درباره دلایل shutdown نوشتم را در نمودارهای pie می‌توانید مشاهده کنید، برای هر سه سوال این نمودارها را رسم شده‌اند. مراجعه کنید به:

[https://github.com/torvalds/linux/blob/4dd58158254c8a027f2bf5060b72ef64cfa3b9d/include/linux/mm\\_types.h#L740](https://github.com/torvalds/linux/blob/4dd58158254c8a027f2bf5060b72ef64cfa3b9d/include/linux/mm_types.h#L740)

متأسفانه بنده خیلی اطلاعات دقیقی درباره تعریف عبارات پیش رو نیافتم، بنابراین با چاشنی عدم اطمینان مطالعه کنید.

**TLB Remote Shutdown:** زمانی رخ می‌دهد که یک هسته‌ی دیگر پردازنده درخواست shutdown را داده باشد.

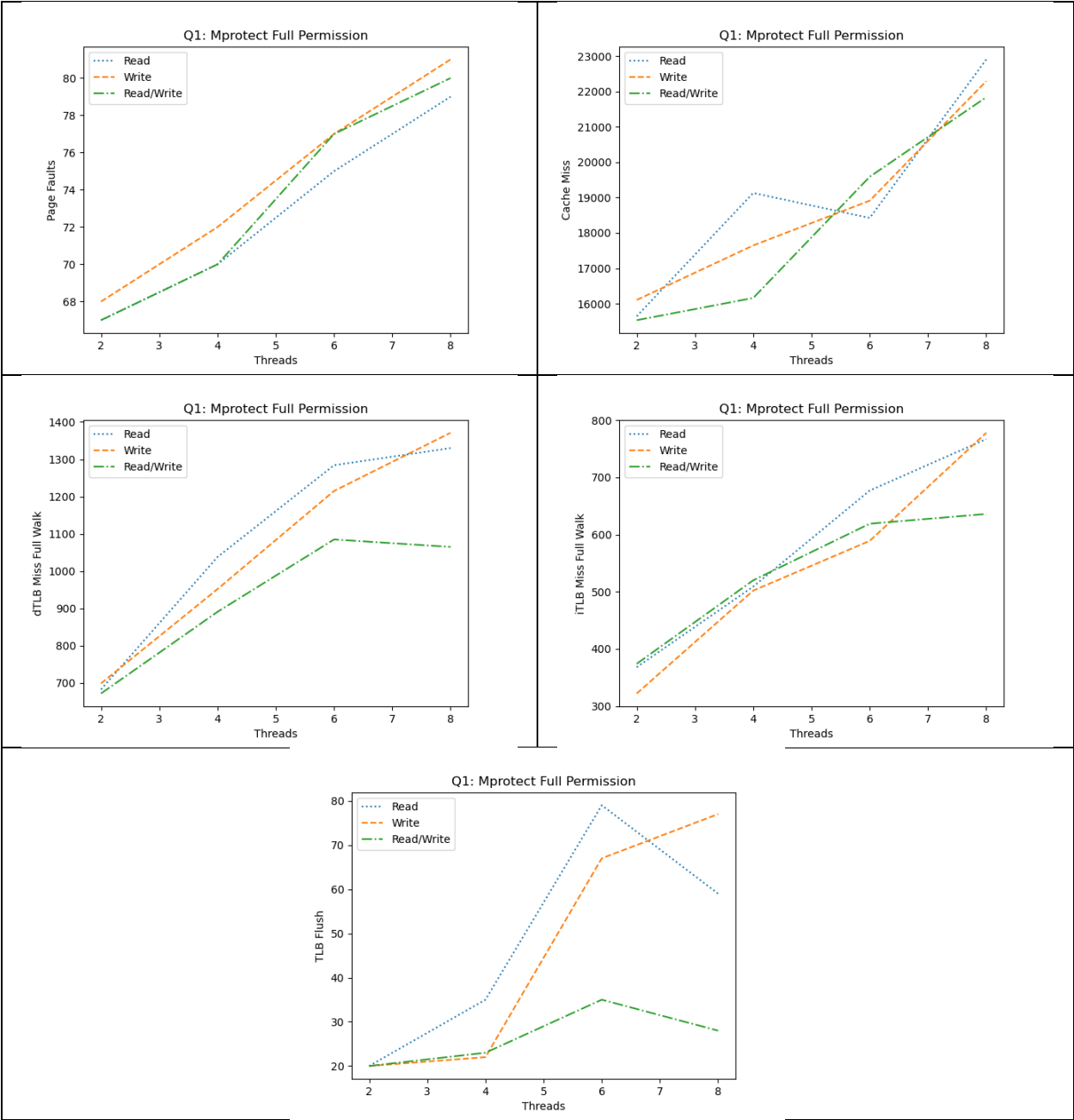
**TLB Remote Send IPI:** فکر می‌کنم این سیگنال و بالایی باید علت و معلول یک دیگر باشد، یعنی پس از ارسال ipi باید Remote shutdown صورت بگیرد. اما جلوتر خواهیم دید که تعداد آنها برابر نیستند. بنابراین مطمئن نیستیم تفاوت کجاست.

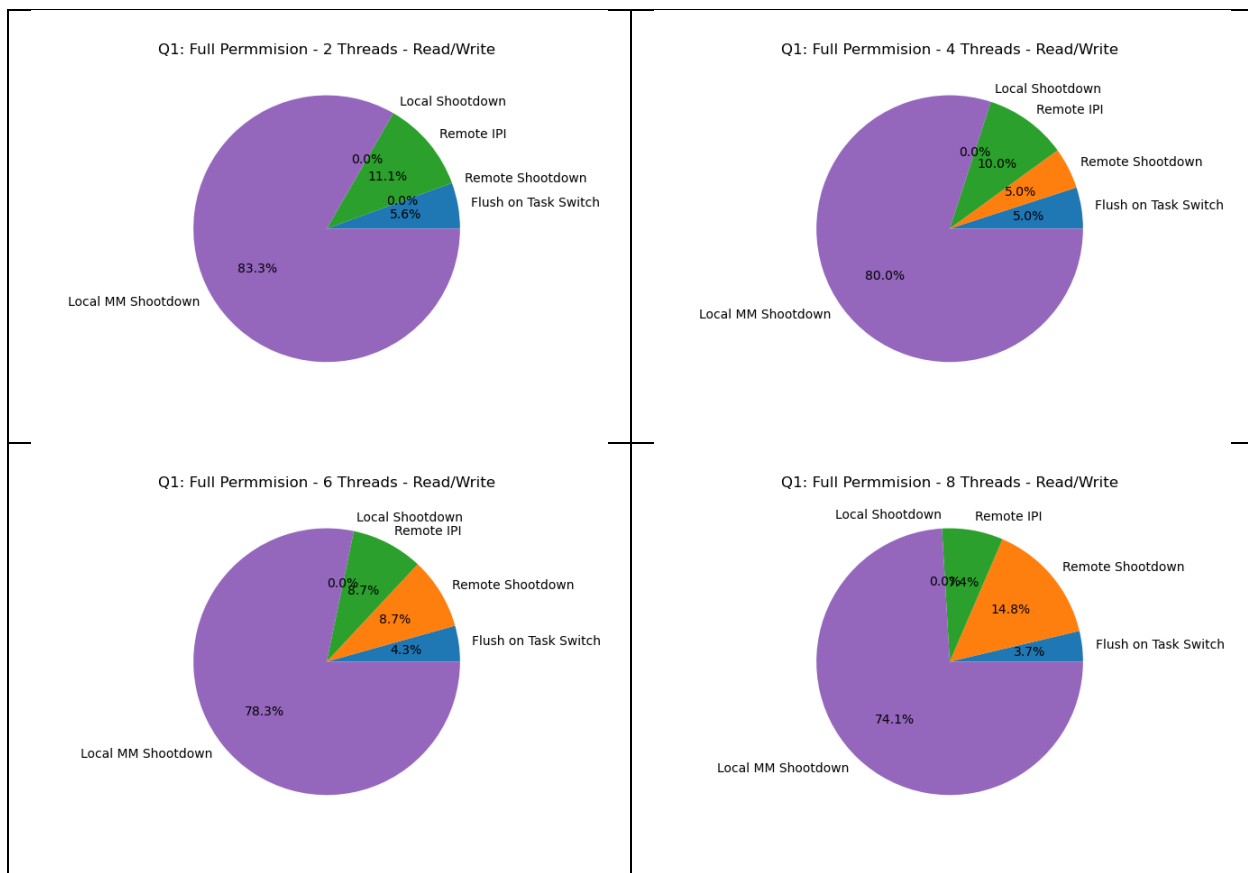
**TLB Local MM Shutdown:** اگر ساختار صفحه توسط هسته فعلی تغییر کند آنگاه این نوع shutdown رخ می‌دهد.

**TLB Local Shutdown:** هنگامی که صفحات نگاشت معکوس شوند یعنی از حافظه اصلی خارج شوند.

**TLB Flush On Task Switch:** از اسم مشخص است که هنگام تعویض متن رخ می‌دهد، جلوتر می‌بینیم که همه‌ی برنامه‌ها حتما این نوع shutdown را تجربه می‌کنند.

پرسش عملی یکم – با دسترسی PROT\_READ, PROT\_WRITE, PROT\_EXEC

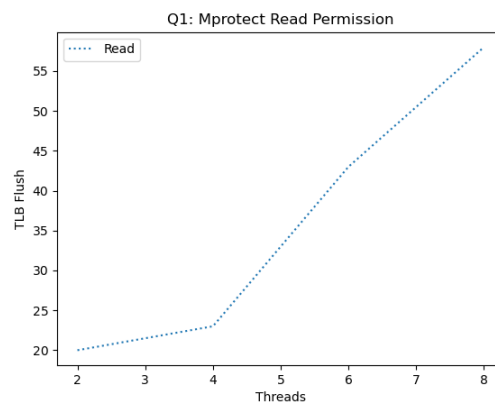
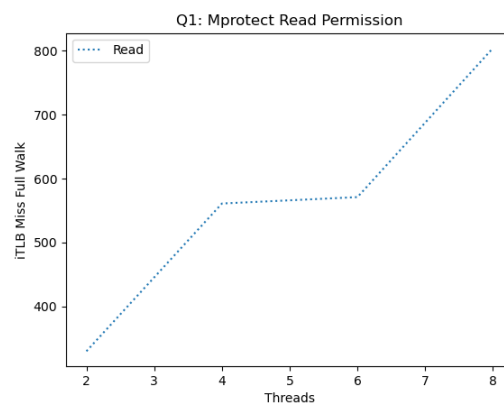
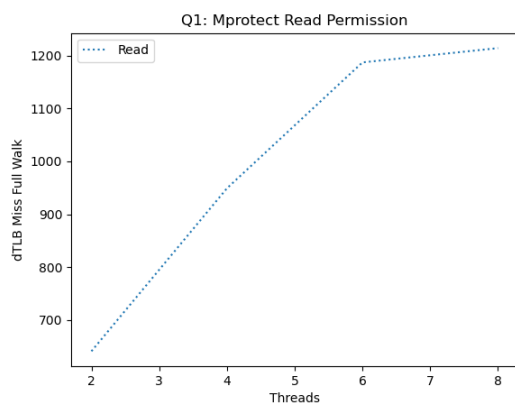
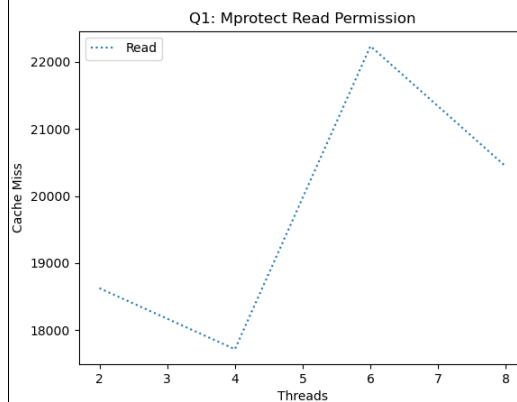
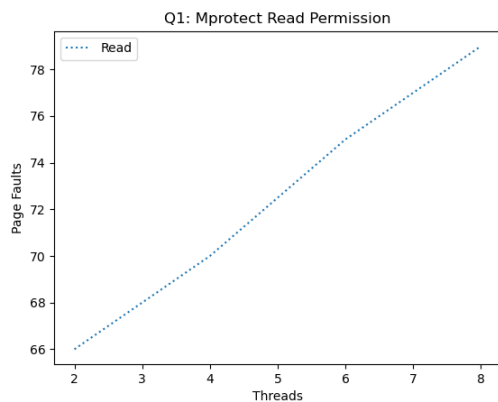




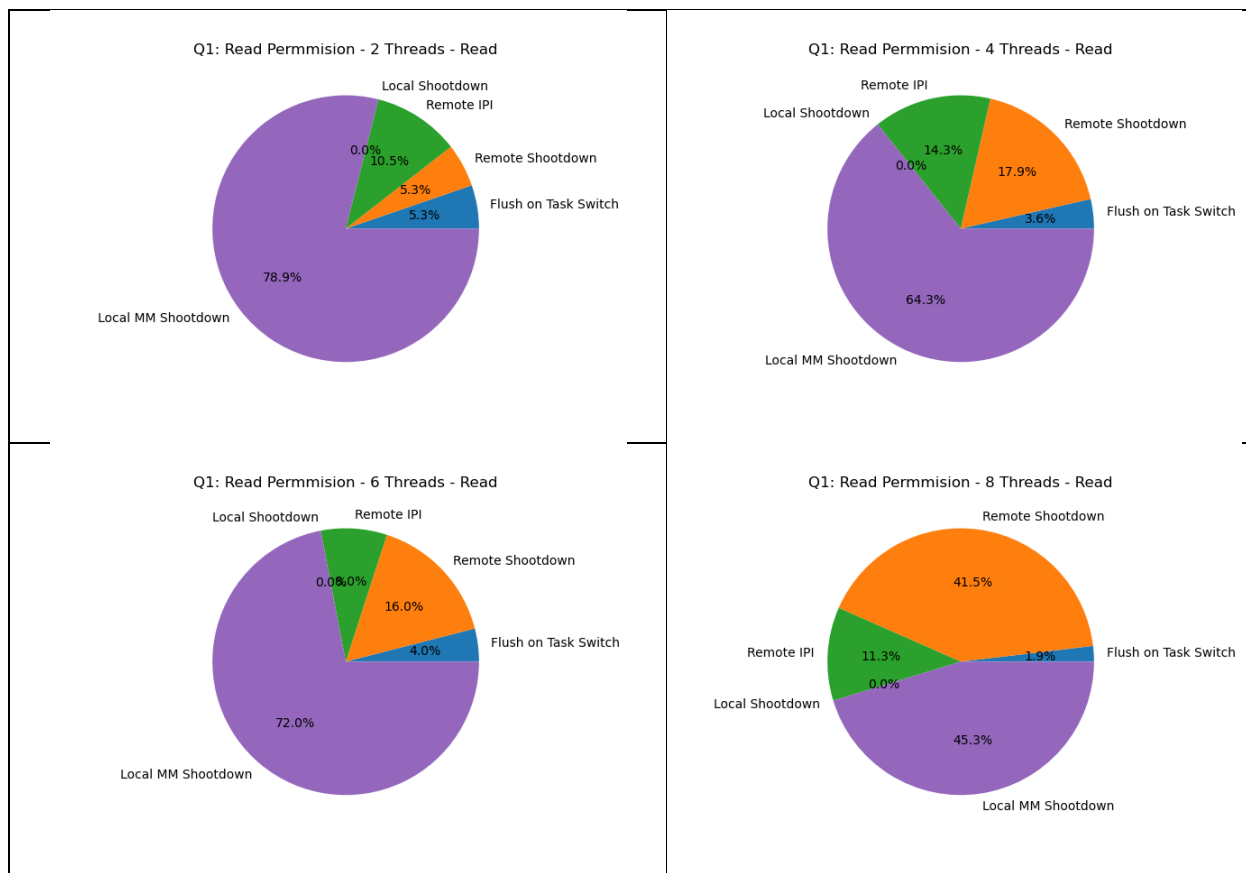
در نمودارهای خطی بالا تعداد تعداد ترد 2 در حالت Read/Write یعنی یک ترد میخواند و دیگری می نویسد، 8 ترد یعنی 4 – 4 و بقیه نیز به همین ترتیب. در حالت فقط Read و فقط Write تمام تردها یک عملیات را انجام می دهند. شیوی اجرا بدین صورت است یک بار تمام تردها در حافظه مورد نظر عملیات خود را انجام می دهد و پشت barrier منتظر می ماند تا همه عملیات را به پایان رسانیده باشند. سپس با استفاده از دستور mprotect سطح دسترسی تغییر می کند و دوباره تمام تردها عملیات خود را انجام می دهند. در این حالت من کامل دسترسی ها را دادم تا بتوانم تمام عملیات ها را انجام دهم. اگر روند کلی نمودارهای خطی را نگاه کنیم بنظر می رسد که با افزایش تردها میزان page fault ها و miss ها و shutdown افزایش پیدا می کند، البته گاهی ممکن است کاهش نیز رخ دهد، من چندین بار مثلا برای تعداد 6 ترد را اجرا کردم اما واقعا الگوی مشخصی وجود ندارد. یعنی نمی توان نتیجه گیری کرد که مثلا با فلان تعداد ترد ممکن است میزان miss و shutdown کاهش پیدا کند. تمامی حالات دقیقا دارای یک عدد TLB flush on task switch هستند، که خب طبیعی به نظر می رسد. تمامی page fault ها از نوع minor page fault هستند. یعنی اینکه داده ی مورد نظر در حافظه اصلی قرار دارد اما رکورد ترجمه آن نیست، البته دلیل دیگر نیز می تواند داشته باشد مانند سطح دسترسی [Architectural and Operating System Support for Virtual memory – page 29].

اگر برنامه بالا را بدون چندنخی و مثلا فقط با یک نخ انجام دهیم و آنگاه میزان shutdown را اندازی گیری کنیم متوجه می شویم که تعداد ipi ها و remote ها به صفر می رسد که البته منطقی است. من این را با record تست کردم ولی در نمودارها نیاوردم.

## پرسش عملی یکم - با دسترسی PROT\_READ



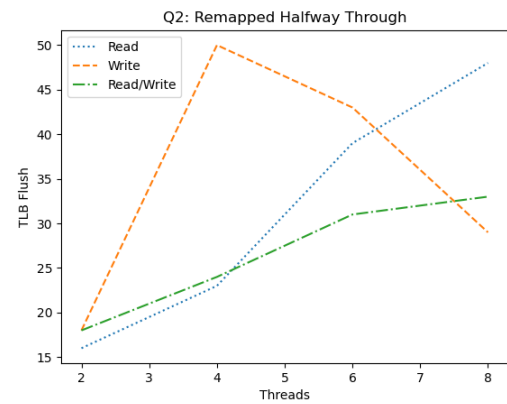
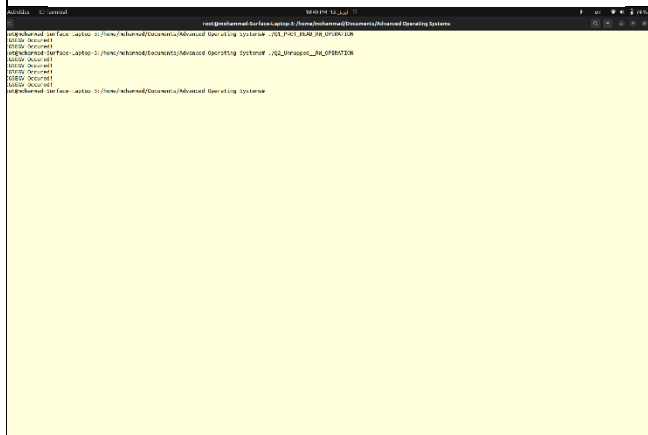
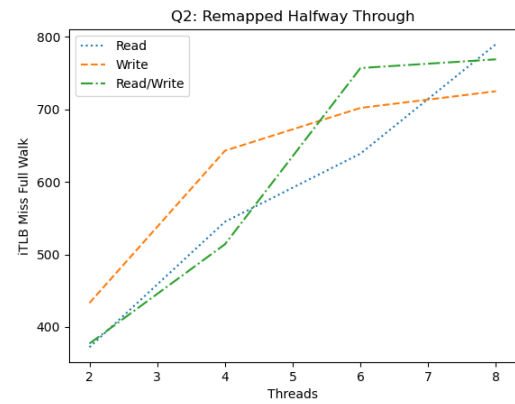
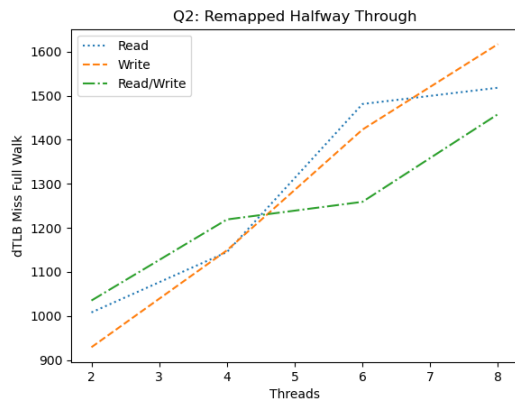
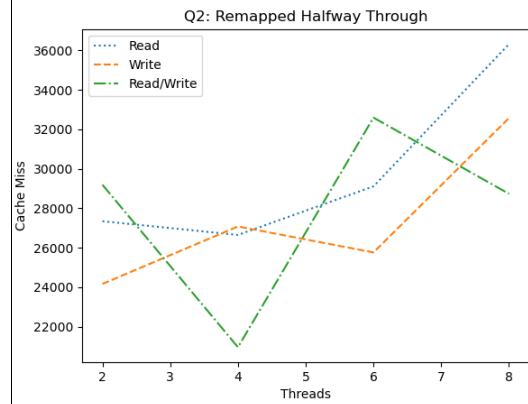
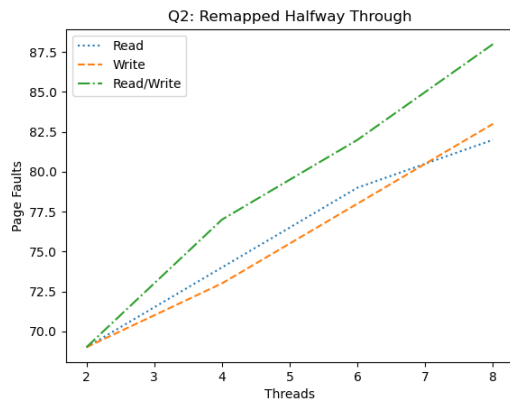
```
root@kali: ~# cat /etc/passwd | grep root
root:x:0:0:root:/root:/bin/bash
root@kali: ~# cat /etc/passwd | grep root
root:x:0:0:root:/root:/bin/bash
```



در این حالت ترتیب اجرا مانند قبل است اما سطح دسترسی را PROT\_READ تعیین کردم، بنابراین فقط امکان خواندن وجود دارد. اگر پس از تغییر دسترسی بخواهیم عملیات نوشتن انجام دهیم سیگنال SIGSEGV صادر شده و خطای Segmentation fault را در میان برنامه و هنگام تلاش برای نوشتن دریافت خواهیم کرد. با استفاده از یک signal handler می‌توان آنرا دریافت کرده و عبارت SIGSEGV OCCURED را در صورت وقوع در خروجی خطای استاندارد نمایش داد. اسکرین‌شات مربوطه در جدول بالا قابل مشاهده است.

محسوس‌ترین تغییر را می‌توان در نمودار pie مشاهده کرد، اگر توجه کنید تعداد Remote Shutdown نسبت به حالت قبلی که عملاً سطح دسترسی تغییری نمی‌کرد به شدت افزایش یافته، دلیلش محدود شدن سطح دسترسی توسط تردی به نام thread\_protector (مراجعه شود به فایل Q2.c) که درخواست Shutdown را به سایر هسته‌ها ارسال می‌کند. (شاید هم تصادفی بوده).

در این حالت نیز تمامی fault page ها minor هستند. همانطور که پیشتر گفتم این نوع نقص به دلیل عدم دسترسی به آدرس است ولی خود داده در حافظه اصلی قرار دارد.





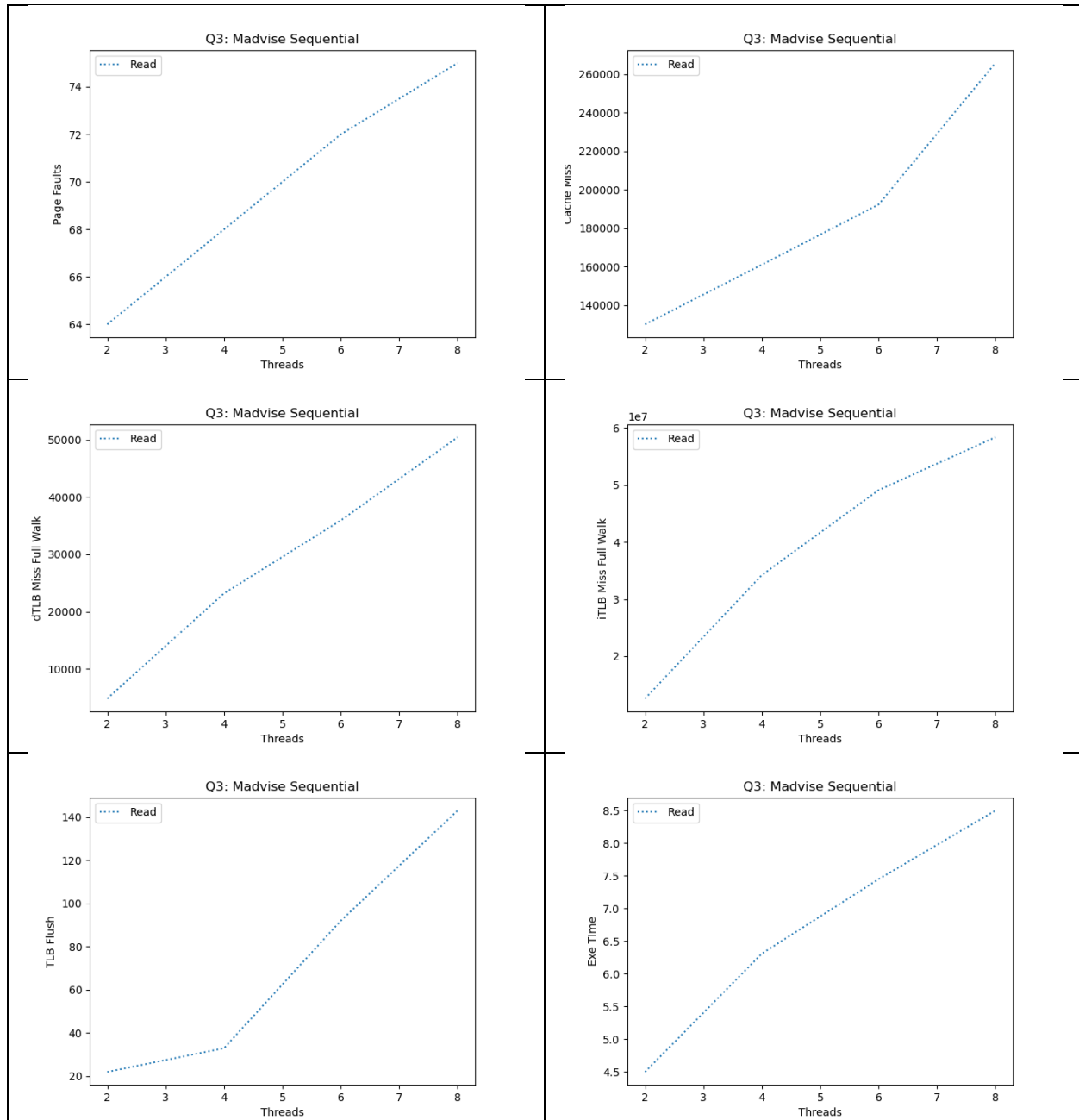
ترتیب اجرا مانند حالات پیشین و با استفاده از barrier صورت گرفته. یعنی یک بار تمامی تردها در نگاشت اول و بار دیگر همگی در نگاشت ثانویه عملیات انجام می‌دهند. اگر در این میان عمل نگاشت معکوس صورت بگیرد خطای Segmentation دریافت خواهیم کرد، اسکرین‌شات در جدول فوق این را نشان می‌دهد. دلیلش نیز بدیهی است زیرا به عملاً به قسمتی تلاش دسترسی داریم که دیگر در اختیار ما نیست.

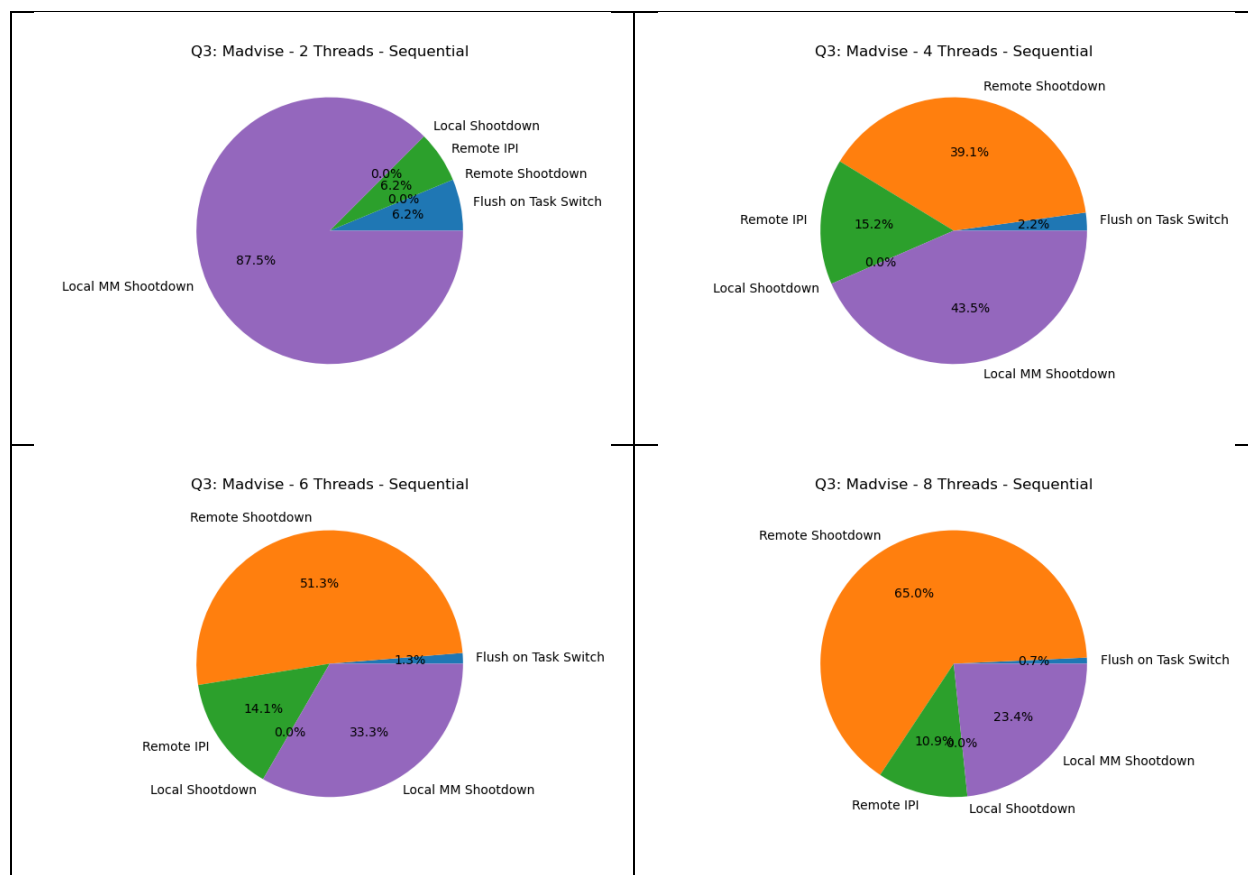
بنظر مهم‌ترین نکته در این سوال در Page fault هاست، این تنها حالتی است Major Page Fault دریافت می‌کنیم. تعداد دقیق بدین شرح است:

Major PF	2 Threads	4 Threads	6 Threads	8 Threads
Read	3	4	4	5
Write	3	4	4	6
R / W	4	7	7	9

Major Page Fault اینطور تعریف می‌شود: تلاش برای دسترسی به داده‌ای که در حافظه اصلی وجود ندارد و باید از حافظه جانبی به حافظه اصلی آورده شود. از آنجایی که در این سوال حافظه جانبی را نگاشت می‌کنیم به نقص‌ها بنظر کاملاً منطقی می‌رسند. افزایش تعداد نقص‌ها احتمالاً به این دلیل است که پس از میزان مشخصی نوشتن سیستم‌عامل محتویات نگاشت شده حافظه اصلی را به جانبی منتقل کرده و دوباره به اصلی بر می‌گرداند.







شاید بهتر بود به صورت دیگر از `madvise` استفاده می‌شد، مثلاً آنرا با آپشن `Sequential` اجرا کرد سپس به صورت `Scan` و `Sequential` خواند یا نوشت تا تفاوت آن با حالت عادی مشخص شود.

من این دستور را با آپشن `Sequential` اجرا کردم، این آپشن باعث می‌شود که پس از یک بار خواندن یا نوشتن داده‌ها به سرعت از `Cache` و یا `TLB` اخراج شوند و آنهایی که جلوتر هستند وارد شوند. دلیلش نیز واضح است، چون یک بار از اول تا آخر یک چیزی را می‌خواهیم خوانیم یا بنویسیم.

این امر در نمودار `pie` مشخص است، به تعداد زیاد `Remote Shootdown` توجه کنید که با بالا رفتن تعداد تردها که هر کدام دارند `madvise` را فراخوانی می‌کنند به شدت افزایش پیدا می‌کند.

توضیحاتی درباره تمرینات عملی: برای دقیق تر بودن داده‌های استخراجی پیش از اجرای هر دستور **perf** و اجرای برنامه‌ها، حافظه‌ی نهان صفحه و داده را خالی کردم. همچنین دستورات **stat** را با فراوانی 90 بار اجرا کردم تا در حد امکان از واریانس بالا و ناهنجاری در داده‌ها جلوگیری شود. فقط پرسش سوم عملی را به دلیل طولانی بودن زمان اجرا مجبور شدم با فراوانی 1 اجرا کنم.

برای بدست آوردن دلایل **tlb shutdown** از آپشن **record** بهره بردم سپس با استفاده از **script** آن را به فایل قابل خواندن تبدیل کردم. دستورات مورد استفاده به شرح زیر می‌باشند:

```
echo 3 > /proc/sys/vm/drop_caches; perf stat -r 90 -e faults,major-faults,minor-faults,tlb:tlb_flush,dtlb_load_misses.walk_completed,itlb_misses.walk_completed,cache-misses,dTLB-load-misses,iTLB-load-misses
```

```
echo 3 > /proc/sys/vm/drop_caches; perf record -e tlb:tlb_flush
```

```
perf script -i perf.data > foolan.txt
```

اکثر اسکرین‌شات‌های گرفته شده در فایل پیوست موجود هستند و می‌توانید مشاهده کنید اما از آنجا که تعداد آنها بسیار زیاد بوده از آوردن آنها به داخل گزارش خودداری نمودم، توجه کنید که فقط برای نمودارهای یک قسمت چندین بار برنامه با تنظیمات مختلف کامپایل و با دو حالت مختلف **stat** و **record** هر کدام از آنها اجرا شده‌اند. این کار برای تمامی قسمت‌ها انجام پذیرفته است. در مجموع حالات شاید بیشتر از 80 بار. بنابراین تعداد محدودی اسکرین شات احتمالا داخل پیوست موجود نباشد.