



پانخ سوال (۱)

$$\frac{256 \times 2^{10} \times 50 \text{ ns}}{40 \text{ ms}} \times 100 = 32/768 \text{ درصد}$$

پانخ سوال (۲)

نرخ فقدان زمانی کمینه است که دستورهایی که وارد حافظه نهان شده‌اند در طول اجرای برنامه از حافظه نهان خارج نشوند. برای این امر باید هیچیک از دو بلوک دستور، به یک بلوک از حافظه نهان نگاشت نشوند، یعنی باید تعداد کل دستورهای برنامه از اندازه حافظه نهان کوچکتر باشد.

اندازه حافظه نهان یک کیلوبایت یعنی ۲۵۶ کلمه می‌باشد. هر دستور نیز یک کلمه است پس ظرفیت حافظه نهان ۲۵۶ دستور است. پس:

$$I_{max} = 256 - 3 = 253$$

توجه: اگر اولین دستور از خانه صفر حافظه شروع نشود باید فاصله اولین دستور برنامه تا ابتدای بلاکی که در آن قرار می‌گرفت از عدد بالا کم شود. مثلاً اگر آدرس اولین دستور در خانه ۲۰ حافظه بود، از ابتدای بلوک متناظرش که بلوک صفر می‌باشد به اندازه ۱۲ بایت یعنی $\frac{20}{4} = 5$ کلمه فاصله دارد لذا $I_{max} = 256 - 3 - 5 = 248$ و اگر اولین دستور در خانه ۱۵۲ از حافظه اصلی قرار داشت، از

$$I_{max} = 256 - 3 - 6 = 247 \text{ در نتیجه}$$

پانخ سوال (۳)

برای محاسبه زمان اجرا کافی است تعداد سیکل‌های لازم برای اجرای برنامه را بدون تاثیر سلسله مراتب حافظه با تعداد سیکل‌های تعلق به خاطر فقدان جمع کرده و در پرپود کلاک ضرب کنیم.

برای فهم بهتر سوال، در سه حالت زمان اجرا را حساب می‌کنیم:

$$I = I_{max} \text{ حالت اول}$$

تمامی دستورها در حافظه نهان جای می‌گیرند. ابتدا از تاثیر سلسله مراتب حافظه صرفنظر کنیم. دستور اول که مقدار ۱۰۰ را در $t1$ قرار می‌دهد، یک بار اجرا می‌شود اما بقیه دستورها به تعداد دفعات اجرای حلقه یعنی ۱۰۰ بار اجرا می‌شوند. لذا تعداد سیکل‌ها برابر است

$$\text{با } 255 \times 100 + 1$$



برای محاسبه تعداد سیکلهایی که صرف انتقال به حافظه نهان شده است کافی است تعداد فقدان‌ها را حساب کرده و در جریمه فقدان ضرب کنیم. جریمه فقدان طبق سوال ۱۰۰۰ سیکل است. چون اندازه برنامه با اندازه حافظه نهان برابر است برای هر بلوک از دستورها یک بار فقدان داریم. کلاً ۲۵۶ دستور داریم و بلوک‌ها ۳۲ دستوری هستند پس ۸ بلوک داریم. پس تعداد فقدان‌ها ۸ است. لذا تعداد سیکل‌های جریمه برابر با ۸×۱۰۰۰ است. پس زمان اجرا برابر است با:

$$CPU\ Time = (1 + 255 \times 100 + 1000 \times 8) \times 1ns = 33/50.1\mu s$$

(آ) حالت دوم، $I = I_{max} + 1$ بدون استفاده از **Loop Fusion**

تعداد سیکل اجرا بدون در نظر گرفتن تاثیر سلسله مراتب حافظه همان تعدادی است که در قسمت قبل محاسبه شد. اما برای محاسبه تعداد فقدان‌ها، دستورها ۲۵۷ دستور هستند. در نتیجه ۹ بلاک را اشغال خواهند کرد. اما حافظه نهان دارای ۸ بلاک است. چون دستورها پشت سر هم هستند و نگاشت مستقیم است، پس بلاک صفر و ۸ دستور به یک نقطه از حافظه نهان نگاشت می‌شوند. پس بلاک‌های ۱ تا ۷ از دستور تنها یکبار فقدان خواهند داشت، اما بلاک‌های صفر و ۸ به تعداد دفعات اجرای حلقه فقدان خواهند داشت. پس تعداد فقدان‌ها $2 \times 100 + 7$ و جریمه آن $1000 \times (2 \times 100 + 7)$ می‌باشد. پس زمان اجرا برای این حالت برابر است با:

$$CPU\ Time = (1 + 255 \times 100 + 1000 \times (7 + 100 \times 2)) \times 1ns = 232/50.1\mu s$$

(ب) حالت سوم، $I = I_{max} + 1$ با استفاده از **Loop Fusion** و شکستن به دو حلقه با $I - 1$ دستور و یک دستور

اگر از تکنیک **Loop Fusion** استفاده کنیم باید یک حلقه را به دو حلقه با همان محدوده اندیسی بشکنیم. لذا ساختار دو حلقه مشابه و به صورت زیر است:

<pre>addi \$t1,\$t1,100 L: addi \$t1,\$t1,-1 253 Instructions here bne \$t1,\$0,L</pre>
<pre>addi \$t1,\$t1,100 L: addi \$t1,\$t1,-1 1 Instruction here bne \$t1,\$0,L</pre>

در واقع کد بالا همانند دو برنامه مجزا عمل می‌کند. قسمت اول همان حالت اول است که در بالا بررسی شد! قسمت دوم نیز مشابه است و به صورت زیر حساب می‌شود:

تعداد سیکل‌های ساعت بدون در نظر گرفتن سلسله مراتب حافظه برابر است با $100 \times 3 + 1$. تنها یک بار فقدان داریم. پس سیکل جریمه 1000×1 است. پس تعداد سیکل لازم برای قسمت دوم $1000 \times 1 + 100 \times 3 + 1 = 1301$ می‌باشد. زمان اجرای این حالت برابر است با:



$$CPU\ Time = 33/50.1\ \mu s + 130.1 \times 1\ ns = 34/80.2\ \mu s$$

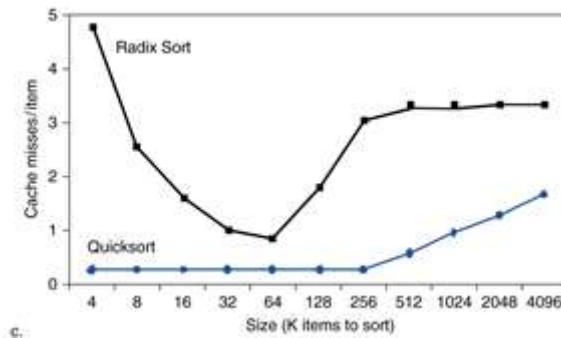
اگر مقدار دو حالت دوم و سوم را مقایسه کنید، مشاهده می‌کنید که در حالت دوم زمان اجرا تقریباً ۲۳۳ میکروثانیه است در حالیکه در حالت سوم حدود ۳۵ میکروثانیه است یعنی با وجود اینکه تعداد دستورهای حالت سوم بیشتر است و برخی دستورها به جای ۱۰۰ بار، ۲۰۰ بار اجرا می‌شوند، اما بازهم زمان اجرای حالت سوم کمتر است. نسبت زمان اجرای حالت دوم و سوم حدود ۷ برابر است.

این مثال تاثیر آشنایی با سخت‌افزار در کارایی نرم‌افزارها را نشان می‌دهد. (و این یکی از دلایل نیاز به یادگیری معماری کامپیوتر برای دانشجویان گرایش نرم‌افزار است)

پانخ سوال (۴)

زمان اجرا تحت تاثیر سه عامل فرکانس، CPI و تعداد دستورها می‌باشد. در یک پردازنده ثابت، فرکانس برای هر دو الگوریتم یکسان است. اگر چه مرتب‌سازی مبنایی دارای تعداد دستورهای کمتری است، اما دارای CPI میانگین بیشتری است و بیشتر بودن CPI به حدی است که زمان اجرای آنرا بیشتر می‌کند.

دلیل بیشتر بودن CPI، وجود فقدان‌های بیشتر در اجرای مرتب‌سازی مبنایی است. قسمت آخر از همان شکل ۵-۱۸ این مطلب را نشان می‌دهد:



اگر فقدان اتفاق افتد، باید تعدادی سیکل تلف شود تا داده مورد نظر از حافظه اصلی به حافظه نهان انتقال یابد و این به معنی افزایش CPI برای اجرای یک دستور است. در واقع عاملی که روشهای کلاسیک و استاندارد تحلیل الگوریتم‌ها از آن چشم‌پوشی می‌کنند، **تاثیر سلسله مراتب حافظه در زمان اجرای یک الگوریتم** است! (این نیز دلیل دیگری برای نیاز به یادگیری معماری کامپیوتر برای دانشجویان گرایش نرم‌افزار است)



پانچ‌سوال (۵)

حافظه مجازی هر نگاشتی که داشته باشد، هر سطر آن دارای ساختار زیر است (طبق صورت سوال از بیت‌های LRU چشم‌پوشی شد):

Valid Bit	Tag	Block Data
-----------	-----	------------

در این سوال، حافظه نهان یک حافظه ۱۶ سطری ۱۰۲۴ کلمه‌ای است. یعنی حافظه نهان این سوال ۱۶ سطر همانند شکل بالا دارد که بخش Block Data آن ۱۰۲۴ کلمه‌ای است. پس اگر تعداد بیت‌های لازم برای Tag را T و اندازه هر کلمه را w فرض کنیم، اندازه حافظه نهان برابر است با:

$$\text{اندازه حافظه نهان} = 16 \times (1 + T + 1024w)$$

در حالت کلی آدرس یک خانه از حافظه را می‌توان به فیلدهای زیر تقسیم کرد:

Tag	Set	Block Offset	Byte Offset
-----	-----	--------------	-------------

در این سوال حافظه اصلی ما دارای واحد آدرس‌پذیر کلمه است. پس Byte Offset دارای صفر بایت خواهد بود. تعداد کلمات داخل یک بلاک ۱۰۲۴ کلمه است پس تعداد بیت‌های Block Offset برابر با $\log_2 1024 = 10$ است. تعداد بیت‌های Set به درجه آزادی بستگی دارد و اگر درجه آزادی a باشد تعداد بیت‌های لازم برای فیلد Set برابر با $\log_2 \frac{16}{a}$ خواهد بود. همچنین حافظه اصلی دارای 2^{20} کلمه است و آدرس‌دهی بر مبنای کلمه است (واحد آدرس‌پذیر کلمه است نه بایت) پس اندازه آدرس برابر با ۲۰ بیت است. پس تعداد بیت‌های لازم برای Tag برابر است با:

$$T = 20 - 0 - 10 - \log_2 \frac{16}{a}$$

آ) نگاشت مستقیم:

نگاشت مستقیم همان نگاشت شبه‌انجمنی با درجه آزادی (درجه انجمنی) یک است. لذا $a = 1$

با جایگذاری در فرمول $T = 6$

$$\text{اندازه حافظه نهان} = 16 \times (1 + 6 + 1024w)$$

آ) نگاشت انجمنی کامل:

نگاشت انجمنی کامل همان نگاشت شبه‌انجمنی است با درجه آزادی ۱۶ (برای این سوال). لذا $a = 16$. پس $T = 10$ در نتیجه:

$$\text{اندازه حافظه نهان} = 16 \times (1 + 10 + 1024w)$$



(آ) نگاشت شبه‌انجمنی با درجه آزادی ۲:

$$a = 2, \text{ پس } T = 7. \text{ لذا}$$

$$\text{اندازه حافظه نهان} = 16 \times (1 + 7 + 1024W)$$

پانچ (۶)

نگاشت مستقیم:

در این نگاشت شماره بلوک در حافظه نهان برابر است با باقیمانده آدرس بلوک در حافظه اصلی بر تعداد بلاک‌ها در حافظه نهان. تعداد بلاک‌ها در حافظه نهان طبق صورت سوال ۴ است. لذا داریم:

شماره بلوک	آدرس بلوک
$\cdot \bmod 4 = 0$	۰
$4 \bmod 4 = 0$	۴
$8 \bmod 4 = 0$	۸
$\cdot \bmod 4 = 0$	۰
$4 \bmod 4 = 0$	۴
$8 \bmod 4 = 0$	۸

چون حافظه نهان در ابتدا خالی است، پس درخواست اول همواره با فقدان روبه‌رو می‌شود. درخواست دوم بلوک ۴ از حافظه اصلی را درخواست می‌کند که باید در بلوک ۰ از حافظه نهان قرار گرفته باشد که وجود ندارد. لذا با فقدان رو به رو می‌شود. درخواست سوم نیز همانند درخواست دوم با فقدان رو به رو می‌شود. درخواست سوم بلوک صفر از حافظه اصلی را می‌خواهد که باید در خانه ۰ از حافظه نهان قرار داشته باشد، اما در این خانه، بلوک ۸ از حافظه اصلی قرار دارد، لذا بازهم فقدان داریم. دو درخواست دیگر نیز با فقدان رو به رو می‌شوند.

نگاشت انجمنی:

در این نگاشت هر بلوک از حافظه اصلی می‌تواند در هر یک از بلوک‌های حافظه نهان باشد یا قرار گیرد. طبق صورت سوال، سیاست جانشینی، LRU است یعنی بلاکی از حافظه نهان جانشین می‌شود اخیراً کمتر استفاده شده باشد. در این سوال برای درخواست‌ها داریم:



۱. بلاک صفر از حافظه اصلی درخواست می شود حافظه نهان خالی است پس فقدان اتفاق می افتد. اکنون باید بلاک درخواستی از حافظه اصلی به حافظه نهان منتقل شود. چون حافظه نهان خالی است، هر جایی از آن می توان این بلاک را قرار داد. ما خانه صفر را انتخاب می کنیم (معمولاً در کامپیوتر شروع از صفر است)
۲. بلاک ۴ از حافظه اصلی درخواست می شود اما در حافظه نهان وجود ندارد پس فقدان داریم. پس از انتقال بلاک درخواستی از حافظه اصلی باید آنرا در بلاک مناسب از حافظه نهان قرار دهیم. بلاک صفر از حافظه نهان پر است اما بلاک ۱ و ۲ خالی است. طبق سیاست LRU بلاک ۴ حافظه اصلی در بلاک ۱ حافظه نهان قرار می گیرد.
۳. بلاک ۸ از حافظه نهان درخواست می شود. باز هم فقدان. بلاک های ۱ و ۲ از حافظه نهان پر است. طبق سیاست LRU این بلاک حتماً باید در بلاک ۲ حافظه نهان قرار گیرد.
۴. اکنون بلاک صفر حافظه اصلی درخواست می شود که در بلاک صفر از حافظه نهان قرار دارد. لذا برخورد اتفاق می افتد.
۵. بلاک ۴ از حافظه اصلی درخواست می شود که در بلاک ۱ از حافظه نهان است. پس برخورد داریم.
۶. بلاک ۸ از حافظه اصلی درخواست می شود که در بلاک ۲ از حافظه نهان جای گرفته است. بنابراین باز هم برخورد داریم.

Hit/Miss	شماره بلوک	آدرس بلوک
Miss	۰	۰
Miss	۱	۴
Miss	۲	۸
Hit	۰	۰
Hit	۱	۴
Hit	۲	۸

نگاشت شبه انجمنی با درجه آزادی ۲:

در این نگاشت ۴ بلاک از حافظه نهان به دو مجموعه دو بلاکی تقسیم می شود. باقیمانده تقسیم آدرس بلاک بر تعداد مجموعه ها شماره مجموعه را مشخص می کند. پس از تعیین شماره مجموعه، هر یک از بلاک های آن مجموعه می تواند حاوی بلاک مورد نظر ما باشد یا بلاک مورد نظر ما در صورت فقدان می تواند با صلاح دید LRU در یکی از آنها جانشین شود. تعداد مجموعه ها برابر است با تعداد بلاک های حافظه نهان تقسیم بر درجه آزادی که برابر با $\frac{4}{2} = 2$ است. تمامی آدرسهای درخواستی دارای باقیمانده ۰ در تقسیم به ۲ هستند. لذا همه درخواست ها به مجموعه صفر از حافظه نهان نگاشت می شوند. پس در یکی از بلاک های صفر یا یک حافظه نهان قرار دارند یا در صورت فقدان در آنها جانشین می شوند. در این سوال برای درخواست ها داریم:

۱. بلاک صفر از حافظه اصلی درخواست می شود حافظه نهان خالی است پس فقدان اتفاق می افتد. اکنون باید بلاک درخواستی از حافظه اصلی به حافظه نهان منتقل شود. چون مجموعه صفر خالی است، هر جایی از آن می توان این بلاک را قرار داد. ما خانه صفر را انتخاب می کنیم (معمولاً در کامپیوتر شروع از صفر است)

۲. بلاک ۴ از حافظه اصلی درخواست می‌شود اما در حافظه نهان وجود ندارد پس فقدان داریم. پس از انتقال بلاک درخواستی از حافظه اصلی باید آنرا در بلاک مناسب از مجموعه صفر از حافظه نهان قرار دهیم. بلاک صفر از حافظه نهان پر است اما بلاک ۱ خالی است. پس حتماً در بلاک ۱ جانشین می‌شود.
۳. بلاک ۸ از حافظه نهان درخواست می‌شود. باید در مجموعه صفر این بلاک را بیابیم. اما بلاک ۸ در این مجموعه قرار ندارد. پس فقدان. هر دو بلاک‌های ۰ و یک مربوط به مجموعه صفر بوده و هر دو پر هستند. طبق LRU باید بلاک صفر جانشین شود چون دیرتر مورد استفاده قرار گرفته است.
۴. اکنون بلاک صفر حافظه اصلی درخواست می‌شود باز هم در مجموعه صفر وجود ندارد. پس فقدان داریم. اما این دفعه بلاک صفر حافظه اصلی در بلاک ۱ حافظه نهان جانشین می‌شود.
۵. بلاک ۴ از حافظه اصلی درخواست می‌شود که در مجموعه صفر وجود ندارد. لذا فقدان اتفاق می‌افتد و طبق LRU باید در بلاک صفر جانشین شود.
۶. بلاک ۸ از حافظه اصلی درخواست می‌شود که در مجموعه صفر وجود ندارد. پس فقدان داریم و طبق LRU بلاک مقصد ۱ است.

Hit/Miss	شماره بلاک	شماره مجموعه	آدرس بلاک
Miss	۰	۰	۰
Miss	۱	۰	۴
Miss	۰	۰	۸
Miss	۱	۰	۰
Miss	۰	۰	۴
Miss	۱	۰	۸

بنابر آنچه گفته شد پاسخ نهایی سوال به صورت زیر است:

2-way set associative		Fully associative		Direct Map		آدرس بلاک مورد نظر
شماره بلاک	Hit or Miss	شماره بلاک	Hit or Miss	شماره بلاک	Hit or Miss	
۰	Miss	۰	Miss	۰	Miss	۰
۱	Miss	۱	Miss	۰	Miss	۴
۰	Miss	۲	Miss	۰	Miss	۸
۱	Miss	۰	Hit	۰	Miss	۰
۰	Miss	۱	Hit	۰	Miss	۴
۱	Miss	۲	Hit	۰	Miss	۸



پانچ سوال (۷)

مقصود از CPI موثر همان CPI کل است. کل CPI لازم برابر است با CPI پایه بعلاوه کلاک‌هایی که به خاطر فقدان در خط لوله تعلیق اتفاق می‌افتد. چون دو سطح حافظه نهان داریم پس باید برای هر یک جداگانه حساب کنیم. پس:

$$\begin{aligned} \text{Effective CPI} &= \text{Total CPI} \\ &= \text{Base CPI} + \text{Primary Stalls per instruction} \\ &\quad + \text{Secondary Stalls per instruction} \end{aligned}$$

$$\text{Effective CPI} = 1 + 0.5 \times 10 + 0.2 \times 400 = 9/5$$

موفق باشید

گروه حل تمرین