



دانشگاه صنعتی شریف
دانشکده مهندسی کامپیوتر
گزارش پروژه معماری کامپیوتر

عنوان:

پروژه ششم: پیاده‌سازی سیاست‌های جایگزینی در حافظه نهان با استفاده از ChampSim

Project Six: Implementation of Cache Replacement
Policies in the ChampSim Simulator

نگارش

محمد مهدی عابدینی ، کسری منتظری هدشی ، محمدپارسا جعفرنژادی ،
سید محمد رضا جوادی

استاد درس

دکتر اسدی

مسئولان پروژه

آقایان صداقت‌گو و غفوری

تیر ۱۴۰۴

فهرست مطالب

| | |
|----|--|
| ۳ | ۱ مقدمه |
| ۳ | ۱-۱ چشم‌انداز گزارش |
| ۴ | ۲-۱ اهداف |
| ۵ | ۲ سیاست‌های جایگزینی |
| ۵ | ۱-۲ شرح مسئله |
| ۶ | ۲-۲ سیاست‌های جایگزینی |
| ۶ | ۱-۲-۲ LRU |
| ۶ | ۲-۲-۲ MRU |
| ۷ | ۳-۲-۲ Random |
| ۷ | ۴-۲-۲ SRRIP |
| ۸ | ۵-۲-۲ LFU |
| ۹ | ۶-۲-۲ Tree-PLRU |
| ۱۰ | ۷-۲-۲ ARC |
| ۱۱ | ۸-۲-۲ LARC |
| ۱۳ | ۳ پیاده‌سازی |
| ۱۳ | ۱-۳ مطالعه مستندات و آشنایی با پیاده‌سازی سیاست‌ها |
| ۱۳ | ۱-۱-۳ Random |
| ۱۴ | ۲-۱-۳ LRU |
| ۱۵ | ۳-۱-۳ SRRIP |
| ۱۶ | ۲-۳ پیاده‌سازی سیاست‌های جدید |
| ۱۶ | ۱-۲-۳ MRU |
| ۱۷ | ۲-۲-۳ Tree-PLRU |

| | |
|----|--------------------------|
| ۱۸ | LFU ۳-۲-۳ |
| ۱۹ | ARC ۴-۲-۳ |
| ۲۲ | LARC ۵-۲-۳ |
| ۲۳ | شبه‌سازی و پارامترها ۴ |
| ۲۴ | تحلیل نتایج ۵ |
| ۲۹ | پیچیدگی‌های پیاده‌سازی ۶ |
| | مراجع ۳۰ |

۱ مقدمه

با افزایش سرعت پردازنده‌ها و محدودیت‌های سرعت حافظه‌های اصلی، فرآیند عملکردی بین این دو به چالشی بزرگ در طراحی سیستم‌های کامپیوتری تبدیل شده است. حافظه‌های نهان به عنوان لایه‌ای میان، در کاهش تأخیر دسترسی به حافظه‌ی اصلی و بهبود عملکرد سیستم نقش کلیدی ایفا می‌کنند. سیاست‌های جایگزین حافظه‌ی نهان تصمیم می‌گیرند که کدام بلوک داده حذف و جایگزین شود. انتخاب بهینه این الگوریتم‌ها تأثیر زیادی بر کارایی سیستم دارد. در این پروژه، هدف پیاده‌سازی و مقایسه چند سیاست جایگزین حافظه‌ی نهان است تا عملکرد آن‌ها در سناریوهای مختلف تحلیل شده و به درک بهتری از مدیریت حافظه‌ی نهان در سیستم‌های کامپیوتری برسیم. ابزار اصلی شبیه‌سازی ما برای بررسی سیاست‌های مختلف، چمپسیم (ChampSim) است که می‌تواند سناریوهای مختلف حافظه نهان و سیاست‌های جایگزین را به صورت دقیق و مفصل شبیه‌سازی کند. با مقایسه نتایج و تحلیل و تفسیر داده‌ها ما می‌توانیم سیاست‌ها را ارزیابی کرده و مفاهیم تئوری درس را در عمل پیاده‌سازی کنیم و نتایج را به چشم ببینیم.

۱-۱ چشم‌انداز گزارش

در این گزارش، قدم به قدم مراحل انجام پروژه را بررسی می‌کنیم. مراحل پروژه عبارتند از:

۱. لینک‌های مرتبط با هر سیاست را مطالعه می‌کنیم تا با نحوه کارکرد و پیاده‌سازی آن‌ها آشنا شویم.

۲. درک خود از نحوه عملکرد این سیاست‌ها را در گزارش آورده و توضیح می‌دهیم.

۳. این سیاست‌ها را در شبیه‌ساز ChampSim پیاده‌سازی می‌کنیم.

۴. شبیه‌ساز را با استفاده از این سیاست‌ها و سیاست‌های Random LRU و SRRIP که توسط خود شبیه‌ساز پیاده‌سازی شده است، پیکربندی کرده و برنامه‌های مختلف را بر روی آن اجرا می‌کنیم.

۵. عملکرد این سیاست‌ها را با پارامترهای مختلف که توسط شبیه‌ساز ارائه می‌شود، مقایسه و تحلیل می‌کنیم.

پس نتایج در این گزارش ما ابتدا به توضیح سیاست‌های جایگزینی می‌پردازیم و بعد از آن نحوه

پیاده سازی را گزارش می دهیم و نهایتاً با شبیه سازی سیاست های مختلف به نتیجه گیری و تحلیل عملکرد هریک می پردازیم

۲-۱ اهداف

در این پروژه ما اهداف زیر را دنبال می کنیم.

- شناخت و آشنایی با شبیه ساز ChampSim و نحوه کار با آن.
- آشنایی با سیاست های مختلف جایگزین در حافظه نهان.
- مقایسه و تحلیل عملکرد سیاست های مختلف حافظه نهان بر عملکرد کل سیستم.

۲ سیاست‌های جایگزینی

در این بخش نتایج مطالعات خود از لینک‌های داده شده را شرح می‌دهیم. به طور خلاصه، در این بخش مطابق گام دوم پروژه انواع سیاست‌ها را نام می‌بریم و آن‌ها را شرح می‌دهیم. قبل از آن که وارد هر کدام به طور خاص شویم، ابتدا بیایید به طور کلی به شرح مسئله به پردازیم. مهم است که بدانیم اصلاً این سیاست‌ها چه معنی دارند و به طور کلی در پاسخ به چه مشکلی این راه‌حل‌ها ارائه شدند.

۱-۲ شرح مسئله

فرکانس کاری CPU بسیار بالاست. در بسیاری از پیاده‌سازی‌ها، مانند پیاده‌سازی چند چرخه‌ای در مرحله حافظه فرصت محدودی برای کار با حافظه وجود دارد و انتظار می‌رود در یک چرخه عملیات انجام شود. با این وجود به دلیل فرکانس کاری بسیار پایین در حافظه‌ها این عمل غیرممکن است. حتی اگر از تکنولوژی‌های جدید برای حافظه استفاده کنیم تا سرعت آن به پردازنده برسد هزینه به شکل سرسام‌آوری بالا می‌رود و طراحی را ناکارآمد می‌کند. اینجا ایده ایجاد واسطه شکل می‌گیرد. حافظه نهان وظیفه دارد تا با ایجاد واسطه بین پردازنده و حافظه مشکل را رفع کند. حافظه نهان خود اندازه محدودی دارد، پس برای قرارگیری بلوک‌های حافظه اصلی ما با محدودیت روبرو هستیم. اینجا چندین سوال مطرح می‌شود:

۱. بلوک در کدام مجموعه قرار گیرد؟

۲. در صورتی که مجموعه متناظر پر بود با کدام بلوک در مجموعه جایگزین شود؟

۳....

ما در این پروژه به سوال دوم می‌پردازیم. هدف این است که با پیدا کردن سیاست مناسب برای موقعیت‌های مختلف بتوانیم نرخ برخورد را حداکثر کنیم. با وجود ارائه الگوریتم‌های مختلف باید توجه داشته باشید که در لایه‌های نزدیک حافظه فرصت محدودی داریم و نباید زمان زیادی صرف خود الگوریتم شود یا هزینه‌های سخت‌افزاری زیادی برایمان بتراشد. با بهینه کردن تمامی این عوامل موثر ما به جستجوی سیاست مطلوب می‌پردازیم. در ادامه چند سیاست که متداولند یا در این پروژه پیاده می‌شوند می‌پردازیم.

۲-۲ سیاست‌های جایگزینی

در این بخش به سیاست‌ها و به مثال‌هایی از آن‌ها می‌پردازیم. هدف این است که درک و توصیف خود از این سیاست‌ها را ارائه دهیم.

روش پیاده‌سازی به طور کلی استفاده از لیست یا داده ساختاری برای آنها است. در هر بار مراجعه در سیاست‌های مختلف براساس اولویت در داده ساختار درج یا حذف انجام می‌شوند و پیچیدگی‌های مختلف بر اساس نوع داده ساختار و الگوریتم جستجو و درج در آن است. علاوه بر این در حافظه نهان به دلیل محدودیت‌های سخت‌افزاری پیچیدگی فضا نیز بسیار مهم است.

LRU ۱-۲-۲

سیاست LRU بلوکی را که بیش‌ترین زمان از آخرین استفاده از آن گذشته است، حذف می‌کند.

این سیاست از Temporal locality بهره می‌برد، به طوری که فرض می‌کند بلوک‌هایی که اخیراً از آن‌ها استفاده شده با احتمال خوبی دوباره نیز استفاده خواهند شد و آن‌هایی که زمان زیادی از استفاده از آن‌ها گذشته است، احتمالاً دیگر در آینده نزدیک مورد استفاده قرار نخواهند گرفت. این سیاست در مواقعی که دسترسی به یک خانه به طور مداوم و در یک مقطع انجام می‌شود ایده‌آل است.

MRU ۲-۲-۲

سیاست MRU بلوکی را که اخیراً استفاده شده است، برای جایگزینی انتخاب می‌کند، برعکس سیاست LRU.

این نوع سیاست از ساختارهای ساده کد بهره می‌برد. مثلاً اگر مراجعه به خانه‌های حافظه در حلقه باشد متغیری که اخیراً استفاده شده یحتمل در تا اجرای بعدی حلقه استفاده نمی‌شود و می‌توان آن را کنار گذاشت. طبیعتاً این نوع دقیقاً خلاف LRU عمل می‌کند و اگر ساختار کد به شکل دسترسی

مداوم به یک خانه باشد نرخ برخورد ما پایین می آید. پس این ساختار به شدت به نوع دسترسی وابسته است.

۳-۲-۲ Random

به طور تصادفی یک بلوک را انتخاب کن و جایگزین کن.

در دنیای ماشین لرنینگ می گویند اگر بین چند فرض، بایاسی نداشته باشیم انتخاب رندوم بهترین امید ریاضی بازدهی را دارد. در حافظه نهان نیز اگر فرض خاصی نکنیم انتخاب رندوم کارآمد است. مهم ترین دلیل کارآمدی این سیاست این است که پیاده سازی آنی بسیار آسان و کم هزینه است چه از نظر سخت افزاری و چه پیچیدگی زمانی.

۴-۲-۲ SRRIP

بلوک هایی که کمترین مقدار ریفرنسیال را دارند جایگزین کن.

در تلاش برای قرار دادن پارامتری برای مشخص شدن ارزش بلوک، روش LRU می گوید تازگی براساس آخرین استفاده و LFU براساس تعداد استفاده. در این روش با معرفی پارامتر دیگری که ارزش گذاری و آپدیت آن متفاوت از دو الگوریتم گفته شده است ما معیار جدیدی برای جایگزینی انتخاب می کنیم. به این معیار مقدار ریفرنسیال گوئیم. الگوریتم به این شکل است: مقادیر ریفرنسیال کران بالا و پایین دارند مثلاً کران بالا مانند ۳. ابتدا مقدار آغازین برای همه خانه ها یکسان است ولی فول نیست مثلاً ۲.

حال دو حالت داریم: در هنگام دسترسی: ناگهان مقدار ران کامل می کنیم. در گذر زمان و دسترسی نداشتن: آن را به مقدار پیش بینی شده ای تغییر می دهیم.

این نوع نگاه باعث ترکیبی از تازگی و تکرار می شود

تصویر زیر یک نوع پیاده سازی خاص در یک مقاله را هم تحت دایاگرام و هم تحت شبه کد نشان


```

// Candidates and their RRPVs are indexed by w in [0, W - 1]
if access is a hit on way w:
    RRPVs[w] = 0
else: // access is a miss
    while (maximum value of RRPVs != 2^M - 1): // 3 if M == 2
        Increment all RRPVs by 1
    // Now there must be at least an RRPV with value 2^M - 1
    Select w as the minimum index s.t. RRPVs[w] == 2^M - 1
    // w-th candidate is selected as the victim
    Replace the old data of candidate w with the new data
    RRPVs[w] = 2^M - 2 // 2 if M == 2

```

The following state-transition diagram illustrates how the RRPVs are managed for 2-bit SRRIP, i.e., when using $M=2$ -bit RRPVs, a common implementation choice.

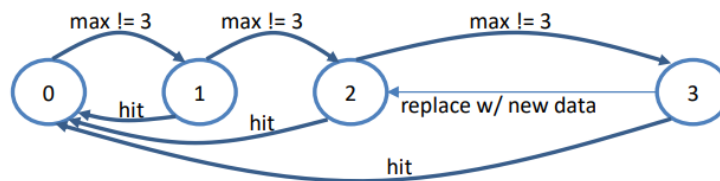


Figure 1. State-transition diagram of an RRPV in 2-bit SRRIP.

LFU ۵-۲-۲

در این سیاست، بلوک‌هایی که کم‌ترین تعداد دسترسی را داشته‌اند، برای جایگزینی انتخاب می‌شوند.

در مبحث Locality معتقدیم گذشته دسترسی‌ها بهترین پیشبینی برای آینده آنها است. در این روش با بهره‌مندی از این اصل ما به نوعی تاریخچه دسترسی به بلوک را نگهداری می‌کنیم و براساس تعداد دسترسی‌ها خانه‌ای که کمترین مراجعه را داشته جایگزین می‌کنیم. در صورت تساوی از LRU استفاده می‌کنیم یعنی خانه‌ای که زودتر وارد شده در صورت برابر بودن تعداد مراجعه اول بیرون می‌رود. این سیاست در الگوهایی که دسترسی به حافظه منظم است و تعداد دسترسی‌ها ناگهانی نیست می‌تواند مفید باشد اما یک عیب بزرگ دارد، اگر توجه کنید بلوک‌هایی که در ابتدا قرار می‌گیرند و به کرات استفاده می‌شوند تا زمان طولانی در حافظه خواهند ماند و بلوک‌های جدیدتر که تازه قرار گرفته‌اند به سرعت خارج می‌شوند. این موضوع به خصوص اگر دسترسی‌ها در یک زمان و پشت سر هم باشد و دیگر به آن خانه رجوع نشود بازدهی بسیار پایینی دارد زیرا جلوی خانه‌هایی که احتمالاً الان استفاده خواهند شد را می‌گیرد و خانه‌هایی که دیگر نیازی به آنها نبود به خاطر نیاز

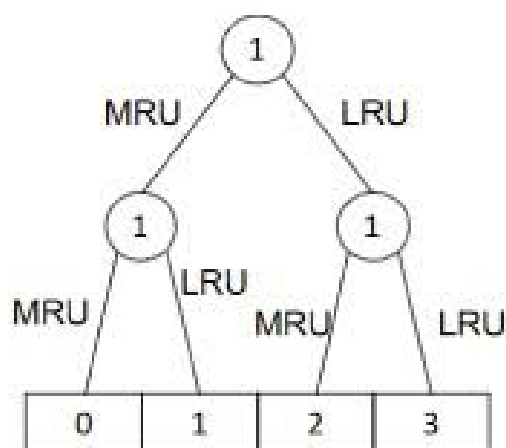
در گذشته جایشان را اشغال می کنند. به نوعی در این سیاست گذشته بسیار مهم است. یک روش پیاده سازی می تواند استفاده از شمارنده ها و آپدیت کردن هر کدام در هربار استفاده از بلوک باشد.

Tree-PLRU ۶-۲-۲

این سیاست از ساختار درختی برای پیگیری وضعیت بلوک ها و تصمیم گیری در مورد بلوک جایگزین بهره می برد و نوعی LRU است

اولین و بدیهی ترین راهی که ممکن است به ذهن هرکسی برسی LRU است. در ادامه ارائه راه حل ها نیز این روش توسعه داده شده یا اشکال متفاوتی برای آن آمده. یک مشکل در روش اصلی این است که تعداد بیت های زیادی صرف نگهداری سن می شود و هزینه زمانی و سخت افزاری زیادی برایمان دارد. یکی از این ایده ها برای بهبود هزینه، استفاده از داده ساختارهای دیگر است. ساختار درختی اجازه می دهد تا هم از تعداد بیت های کمتری استفاده کنیم و هم هزینه جستجوی پایین تری داشته باشیم

یک مثال از نحوه پیاده سازی برای این مورد میزنیم، فرض کنید حافظه نهان ما ۴ سوپیه است. اگر از روش اصلی استفاده کنیم برای دانستن ترتیب آن ها در لیست برای هر بلوک ۲ بیت نیاز داریم، که برای هر مجموعه می شود ۸ بیت. اما حالا بیا ببینیم بجای آن از روش درختی استفاده کنیم. درخت برای هر مجموعه مشابه زیر می شود:



در اینجا ° به معنی این است که زیر شاخه سمت چپ از همه زودتر استفاده شده یا LRU است و ۱ به معنی اینکه زیر شاخه راستی LRU است. به این ترتیب با $\log n$ بیت می توانیم LRU را مشخص کنیم. هنگام جستجو نیز کافی است در این درخت جستجو کرده و از هر مسیر که رفتیم بیت ها را تغییر دهیم. همان طور که می بینید بجای ۸ بیت ۳ بیت استفاده کردیم و هزینه سرشکن تغییر و جستجو و ... هم بسیار کمتر از حالت اصلی است. در نهایت مانند سیاست اصلی بلوکی که کمترین استفاده را اخیراً داشته پس از جستجو جایگزین می شود

ARC ۷-۲-۲

سیاست LRU تمرکز بر تازگی دارد و سیاست LFU تمرکز بر تکرار. هردو مزایای خود را دارند. سیاست ARC با ترکیب این دو و مدیریت هوشمند حافظه نهان برای هر بلوک سیاست بهینه را انتخاب می کند

در این سیاست ما ابتدا کش را به دو قسمت تقسیم می کنیم. بخش اول بخشی است که تمرکز بر تازگی دارد و سیاست LRU در آن اجرا می شود. بخش دوم برای اجرای سیاست تکرار است و سیاست LFU اجرا می شود. اندازه این بخش ها ثابت نخواهد بود ولی در آغاز برایشان مرزی تعیین می کنیم که بعداً طی عملیات ها جابجا می شود.

• T_1 (recently used)

– آیتم هایی که فقط یک بار اخیراً استفاده شده اند.

• T_2 (frequently used)

– آیتم هایی که بیش از یک بار استفاده شده اند.

• B_1 (ghost list for T_1)

– آیتم هایی که از T_1 حذف شده اند، اما هنوز در حافظه کش نیستند؛ فقط نام آن ها نگه داشته می شود (نه داده شان).

• B_2 (ghost list for T_2)

- مشابه B_1 ، اما برای آیتم‌های حذف شده از T_2 .

T_1 و B_1 با همدیگر بخش اول، و T_2 و B_2 بخش دوم را می‌سازند. همچنین در الگوریتم ARC یک پارامتر تطبیقی به نام p دارد که نسبت بین T_1 و T_2 را تعیین می‌کند. اگر آیتم‌های B_1 زیاد hit بخورند (یعنی آیتم‌هایی که فقط یک بار استفاده شده‌اند دوباره درخواست شوند)، مقدار p را افزایش می‌دهد. با اینکار فضای بیشتری به T_1 اختصاص می‌یابد که باعث می‌شود توجه به recency افزایش یابد. به طرز مشابه با hit شدن آیتم‌های B_2 ، مقدار p کاهش می‌یابد تا اندازه ماکسیمم T_2 زیاد شود و باعث افزایش توجه به Frequency خواهد شد بدین ترتیب به صورت کش را به صورت پویا بین دو لیست اصلی تقسیم می‌کند.

۸-۲-۲ LARC

نسخه تغییر یافته ARC است. تفاوت این است که به محض اتفاق افتادن Cache miss حافظه نهان را به روزرسانی نمی‌کند.

در تصویر زیر از مقاله اصلی کلیت پیاده‌سازی سیاست مشهود است این الگوریتم در حقیقت مشابه ARC دارای دو صف اصلی است که هر دو مطابق FIFO عمل می‌کنند اما یکی از آنها تنها دسترسی‌های ۲ یا بیشتر را نگه می‌دارد. تفاوت این سیاست این است که صف اول (معادل T_1) خارج از حافظه نهان فیزیکی نگهداری می‌شود، بدین معنی که این صف در ساختار سیاست جایگزینی وجود دارد و در آن اشاره‌گرهایی به بلوک‌های حافظه (یعنی بخشی از آدرس فیزیکی) نگهداری می‌شود، اما این بلوک‌ها هنوز درون خود حافظه نهان بارگیری نشده‌اند. در حقیقت این سیاست تنها بلوک‌هایی که بیش از یک بار استفاده بشوند را وارد حافظه نهان می‌کند و در غیر این صورت آنها را خارج از آن نگه می‌دارد.

نام صف حافظه فیزیکی Q و نام صف مجازی Q_r است. حداکثر طول صف مجازی را C_r می‌نامیم که پارامتر مهمی است، چرا که اگر زیاد باشد آنگاه بعد مدتی تعداد زیادی از بلوک‌ها در این صف خواهند بود که باعث می‌شود سیاست مانند LRU عمل کند. از طرفی اگر اندازه آن را کم بگیریم

بلوک‌هایی که کمی قبل‌تر آمده‌اند شانس ورودی به حافظه نهان را ندارند. اما مقدار بهینه با توجه به برنامه‌های مختلف متفاوت است چرا که در هر بخش از هر برنامه فواصل بین استفاده دوباره از یک بلوک متفاوت است. این فاصله Reuse Distance نام دارد و برنامه باید مقدار C_r را با توجه به این تنظیم کند.

برای تنظیم پویای این مقدار از روابطی که در تصویر نوشته شده استفاده می‌شود که به طور کلی بدین ترتیب عمل می‌کند:

- اگر بلوک فعلی درون Q یافت شد، اندازه C_r را کاهش بده. هر چه مقدار قبلی آن بیشتر بود آن را بیشتر کاهش بده.
- در غیر این صورت اندازه C_r را افزایش بده، هر چه مقدار قبلی آن کمتر بود آن را بیشتر افزایش بده.

Algorithm 1: Routine to access block B

```

Initilize:  $C_r = 0.1 * C$ 
if  $B$  is in  $Q$  then
    move  $B$  to the MRU end of  $Q$ 
     $C_r = \max(0.1 * C, C_r - \frac{C}{C - C_r})$ 
    redirect the request to the corresponding SSD block
    return
end if
 $C_r = \min(0.9 * C, C_r + \frac{C}{C_r})$ 
if  $B$  is in  $Q_r$  then
    remove  $B$  from  $Q_r$ 
    insert  $B$  to the MRU end of  $Q$ 
    if  $|Q| > C$  then
        remove the LRU block  $D$  in  $Q$ 
        allocate the SSD block of  $D$  to  $B$ 
    else
        allocate a free SSD block to  $B$ 
    end if
else
    insert  $B$  to the MRU end of  $Q_r$ 
    if  $|Q_r| > C_r$  then
        remove the LRU block of  $Q_r$ 
    end if
end if

```

۳ پیاده‌سازی

در این مرحله ما باید براساس مستندات چمپسیم و سیاست های پیاده سازی شده آماده، سیاست های جدید را پیاده می کردیم. پس این بخش دوزیربخش دارد:

۱. مطالعه مستندات و بررسی پیاده سازی های ازپیش آماده

۲. پیاده سازی سیاست های جدید

۱-۳ مطالعه مستندات و آشنایی با پیاده‌سازی سیاست‌ها

در این مرحله با مطالعه مستندات ، راهنمای پیاده سازی ، سایت مفید برای مقالات و مستندات ، راهنمای خود ریپازیتوری و تمامی مقالات اشاره شده در صورت پروژه با نحوه پیاده سازی سیاست ها در این شبیه ساز آشنا شدیم

ابتدا با بررسی ،LRU Random و SRRIP با نحوه پیاده سازی آشنا شدیم. قبل از آن اشاره کنیم کارکرد تک تک توابع را در داکيومنت های چمپسیم می توانید مشاهده کنید.

Random 1-1-3

الگوریتم این سیاست در بخش پیشین توضیح داده شده است در اینجا به بررسی نحوه پیاده سازی می پردازیم.

هر سیاست برای پیاده سازی نیاز به یک فایل هدر برای پروتوتایپ توابع نیاز دارد و یک فایل cc که توابع را پیاده سازی کند.

```

1 #include "random.h"
2
3 random::random(CACHE* cache) : random(cache, cache->NUM_WAY) {}
4
5 random::random(CACHE* cache, long ways) : replacement(cache), dist(0, ways - 1) {}
6
7 long random::find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const CACHE::BLOCK* current_set, uint64_t ip, uint64_t full_addr,
8 | | | | | | | | | | access_type type)
9 {
10 | return dist(rng);
11 }

```

فایل، هدر

```

1 #ifndef REPLACEMENT_RANDOM_H
2 #define REPLACEMENT_RANDOM_H
3
4 #include <random>
5
6 #include "cache.h"
7 #include "modules.h"
8
9 struct random : public champsim::modules::replacement {
10     std::mt19937_64 rng();
11     std::uniform_int_distribution<long> dist;
12
13     explicit random(CACHE* cache);
14     random(CACHE* cache, long ways);
15
16     // void initialize_replacement();
17     long find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const CACHE::BLOCK* current_set, uint64_t ip, uint64_t full_addr, access_type type);
18     // void update_replacement_state(uint32_t triggering_cpu, long set, long way, uint64_t full_addr, uint64_t ip, uint64_t victim_addr, access_type type, uint8_t
19     // hit);
20     // void replacement_final_stats()
21 };
22
23 #endif

```

فایل پیاده سازی

همانطور که می بینید رندوم یکی را برمی گرداند.

۳-۱-۲ LRU

الگوریتم این سیاست در بخش پیشین توضیح داده شده است.

```

1 #ifndef REPLACEMENT_LRU_H
2 #define REPLACEMENT_LRU_H
3
4 #include <vector>
5
6 #include "cache.h"
7 #include "modules.h"
8
9 class lru : public champsim::modules::replacement
10 {
11     long NUM_WAY;
12     std::vector<uint64_t> last_used_cycles;
13     uint64_t cycle = 0;
14
15 public:
16     explicit lru(CACHE* cache);
17     lru(CACHE* cache, long sets, long ways);
18
19     // void initialize_replacement();
20     long find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
21     champsim::address full_addr, access_type type);
22     void replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
23     access_type type);
24     void update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
25     access_type type, uint8_t hit);
26     // void replacement_final_stats()
27 };
28
29 #endif

```

فایل هدر

```

1 #include "lru.h"
2
3 #include <algorithm>
4 #include <cassert>
5
6 lru::lru(CACHE* cache) : lru(cache, cache->NUM_SET, cache->NUM_WAY) {}
7
8 lru::lru(CACHE* cache, long sets, long ways) : replacement(cache), NUM_WAY(ways), last_used_cycles(static_cast<std::size_t>(sets * ways), 0) {}
9
10 long lru::find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
11 champsim::address full_addr, access_type type)
12 {
13     auto begin = std::next(std::begin(last_used_cycles), set * NUM_WAY);
14     auto end = std::next(begin, NUM_WAY);
15
16     // Find the way whose last use cycle is most distant
17     auto victim = std::min_element(begin, end);
18     assert(begin <= victim);
19     assert(victim < end);
20     return std::distance(begin, victim);
21 }
22
23 void lru::replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
24 access_type type)
25 {
26     // Mark the way as being used on the current cycle
27     last_used_cycles.at((std::size_t)(set * NUM_WAY + way)) = cycle++;
28 }
29
30 void lru::update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip,
31 champsim::address victim_addr, access_type type, uint8_t hit)
32 {
33     // Mark the way as being used on the current cycle
34     if (hit && access_type(type) != access_type::WRITE) // Skip this for writeback hits
35         last_used_cycles.at((std::size_t)(set * NUM_WAY + way)) = cycle++;
36 }

```

فایل پیاده سازی

می بینید که در این پیاده سازی آرایه ای نگهداری می کند که هرچه مقدار آن کمتر باشد یعنی بیشترین مدت از طول عمر آن می گذرد. هرگاه نوبت جایگزینی رسد از اعضای آن مینیمم گرفته و آن را جایگزین می کند. توجه کنید اینجا تابع آپدیت شدن وضعیت نیز پیاده شده که اگر برخورد هم داشت باید سن آن یکی افزوده شود.

۳-۱-۳ SRRIP

الگوریتم این سیاست در بخش پیشین توضیح داده شده است.

```
1 #ifndef REPLACEMENT_SRRIP_H
2 #define REPLACEMENT_SRRIP_H
3
4 #include <stdint>
5 #include <vector>
6
7 #include "cache.h"
8 #include "modules.h"
9
10 struct srrip_set_helper {
11     using rrpv_type = int;
12     static constexpr rrpv_type maxRRPV = 3;
13
14     std::vector<rrpv_type> rrpv_values;
15     rrpv_type& get_rrpv(long way);
16
17     explicit srrip_set_helper(long ways);
18
19     long victim();
20     void update(long way, bool hit);
21 };
22
23 struct srrip : public champsim::modules::replacement {
24
25     std::vector<srrip_set_helper> sets;
26
27     explicit srrip(CACHE* cache);
28     srrip(CACHE* cache, long sets_, long ways_);
29
30     // void initialize_replacement() {}
31     long find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
32                     champsim::address full_addr, access_type type);
33     void update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
34                                 access_type type, uint8_t hit);
35
36     // use this function to print out your own stats at the end of simulation
37     // void replacement_final_stats() {}
38 };
39
40 #endif
```

فایل هدر


```

1 #include "srrip.h"
2
3 #include <algorithm>
4 #include <cassert>
5 #include <unordered_map>
6
7 #include "cache.h"
8
9 srrip::srrip(CACHE* cache) : srrip(cache, cache->NUM_SET, cache->NUM_WAY) {}
10
11 srrip::srrip(CACHE* cache, long sets_, long ways_) : replacement(cache)
12 {
13     std::generate_n(std::back_inserter(sets), sets_, [ways = ways_] { return srrip_set_helper{ways}; });
14 }
15
16 // find replacement victim
17 long srrip::find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
18 | | | | | champsim::address full_addr, access_type type)
19 {
20     return sets.at(static_cast<std::size_t>(set)).victim();
21 }
22
23 // called on every cache hit and cache fill
24 void srrip::update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip,
25 | | | | | champsim::address victim_addr, access_type type, uint8_t hit)
26 {
27     sets.at(static_cast<std::size_t>(set)).update(way, hit);
28 }
29
30 srrip_set_helper::srrip_set_helper(long ways) : rrpv_values(static_cast<std::size_t>(ways), maxRRPV) {}
31
32 auto srrip_set_helper::get_rrpv(long way) -> rrpv_type& { return rrpv_values.at(static_cast<std::size_t>(way)); }
33
34 long srrip_set_helper::victim()
35 {
36     // Find the maximum RRPV
37     auto victim = std::max_element(std::begin(rrpv_values), std::end(rrpv_values));
38
39     // If the maximum element has RRPV less than the maximum, increment everything to the maximum
40     std::transform(std::cbegin(rrpv_values), std::cend(rrpv_values), std::begin(rrpv_values), [diff = maxRRPV - *victim](auto x) { return x + diff; });
41
42     // Return the way index
43     return std::distance(std::begin(rrpv_values), victim);
44 }

```

فایل پیاده سازی

در این پیاده سازی به دلیل پیچیدگی داده ساختار ها و متغیر ها و ... بیشتری استفاده شده و پیاده سازی الگوریتم با استفاده از فرایندهای کمکی انجام گرفته تا آپدیت شدن ها به طور خودکار در هر برخورد یا خطا رفت انجام شود.

۲-۳ پیاده سازی سیاست های جدید

حال وقت پیاده سازی سیاست های جدید است.

MRU ۱-۲-۳

پیاده سازی این سیاست آسان بود زیرا همان LRU را با ماکس گرفتن به MRU می توان تبدیل کرد. پس با تغییر نام و استفاده از تابع ماکسیمم به جای مینم این سیاست را پیاده کردیم.

```

1 #ifndef REPLACEMENT_MRU_H
2 #define REPLACEMENT_MRU_H
3
4 #include <vector>
5
6 #include "cache.h"
7 #include "modules.h"
8
9 class mru : public champsim::modules::replacement
10 {
11     long NUM_WAY;
12     std::vector<uint64_t> last_used_cycles;
13     uint64_t cycle = 0;
14
15 public:
16     explicit mru(CACHE* cache);
17     mru(CACHE* cache, long sets, long ways);
18
19     // void initialize_replacement();
20     long find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
21                     champsim::address full_addr, access_type type);
22     void replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
23                                access_type type);
24     void update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
25                                  access_type type, uint8_t hit);
26     // void replacement_final_stats();
27 };
28
29 #endif

```

فایل هدر

```

1 #include "mru.h"
2
3 #include <algorithm>
4 #include <cassert>
5
6 mru::mru(CACHE* cache) : mru(cache, cache->NUM_SET, cache->NUM_WAY) {}
7
8 mru::mru(CACHE* cache, long sets, long ways) : replacement(cache), NUM_WAY(ways), last_used_cycles(static_cast<size_t>(sets * ways), 0) {}
9
10 long mru::find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
11                      champsim::address full_addr, access_type type)
12 {
13     auto begin = std::next(std::begin(last_used_cycles), set * NUM_WAY);
14     auto end = std::next(begin, NUM_WAY);
15
16     // Find the way whose last use cycle is least distant
17     auto victim = std::max_element(begin, end);
18     assert(begin <= victim);
19     assert(victim < end);
20     return std::distance(begin, victim);
21 }
22
23 void mru::replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
24                                  access_type type)
25 {
26     // Mark the way as being used on the current cycle
27     last_used_cycles.at((std::size_t)(set * NUM_WAY + way)) = cycle++;
28 }
29
30 void mru::update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip,
31                                    champsim::address victim_addr, access_type type, uint8_t hit)
32 {
33     // Mark the way as being used on the current cycle
34     if (hit && access_type != access_type::WRITE) // skip this for writeback hits
35         last_used_cycles.at((std::size_t)(set * NUM_WAY + way)) = cycle++;
36 }

```

فایل پیاده سازی

۲-۲-۳ Tree-PLRU

همان طور که در الگوریتم این یک نسخه تغییر یافته برای LRU است. در فایل هدر علاوه بر توابع، یک وکتور برای ذخیره سازی همه درخت ها وجود دارد. در فایل پیاده سازی می بینید که در ابتدا بر اساس ویژگی های حافظه درخت هارا می سازیم.

در واقع، این رویکرد تلاش می کند تا علاوه بر بهره برداری از منطق، LRU از ساختارهای درختی برای افزایش کارایی، کاهش زمان جستجو، و بهبود تصمیم گیری در فرآیند جایگزینی و مدیریت حافظه بهره مند شود. نتیجه نهایی، الگوریتمی است که نه تنها به صورت خطی بلکه به کمک ساختارهای درخت، عملیات را انجام می دهد و پاسخ های بهینه تری در مواجهه با حالت های مختلف حافظه ارائه می دهد. منطق کلی با LRU یکی است تفاوت در نوع داده ساختار است که اینجا بجای مینم گرفتن از یک لیست یا افزودن مقدار در لیست، به جستجوی درختی یا درج در درخت می پردازیم.

```

1 #ifndef REPLACEMENT_TPLRU_H
2 #define REPLACEMENT_TPLRU_H
3
4 #include <vector>
5
6 #include "cache.h"
7 #include "modules.h"
8
9 class tplru : public champsim::modules::replacement
10 {
11     long NUM_WAY;
12     std::vector<short> all_btrees;
13     /* above is the vector that contains all binary tree
14     sequentially, tree with index n can be found at index n*NUM_WAY */
15
16 public:
17     explicit tplru(CACHE* cache);
18     tplru(CACHE* cache, long sets, long ways);
19
20     // void initialize_replacement();
21     long find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
22                     champsim::address full_addr, access_type type);
23     void replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
24                                access_type type);
25     void update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
26                                  access_type type, uint8_t hit);
27     // void replacement_final_stats()
28 };
29
30 #endif

```

فایل هدر

```

1 #include "tplru.h"
2 #include <iostream>
3 #include <algorithm>
4 #include <cassert>
5
6 tplru::tplru(CACHE* cache) : tplru(cache, cache->NUM_SET, cache->NUM_WAY) {}
7
8 tplru::tplru(CACHE* cache, long sets, long ways) : replacement(cache), NUM_WAY(ways), all_btrees(static_cast<std::size_t>(sets * ways), 0) {}
9
10 long tplru::find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
11                        champsim::address full_addr, access_type type)
12 {
13     // find the tree for current set
14     std::vector<short> btree(all_btrees.begin() + set * NUM_WAY, all_btrees.end() + (set+1) * NUM_WAY);
15     // walk through the tree with given directions to find the block index
16     long ind = 1;
17     while (ind < NUM_WAY) {
18         ind = ind * 2 + btree[ind];
19     }
20     return ind - NUM_WAY;
21 }
22
23 void tplru::replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
24                                   access_type type)
25 {
26     // set a block's parents to point in the opposite direction
27     long ind = way + NUM_WAY;
28     while (ind > 1) {
29         all_btrees[set * NUM_WAY + ind / 2] = (ind % 2) ? 0 : 1;
30         ind = ind / 2;
31     }
32 }
33
34 void tplru::update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip,
35                                     champsim::address victim_addr, access_type type, uint8_t hit)
36 {
37     // same as replacement_cache_fill only with conditions to avoid misses and writebacks
38     if (hit && access_type[type] != access_type::WRITE) {
39         long ind = way + NUM_WAY;
40         while (ind > 1) {
41             all_btrees[set * NUM_WAY + ind / 2] = (ind % 2) ? 0 : 1;
42             ind = ind / 2;
43         }
44     }
45 }

```

فایل پیاده سازی

۳-۲-۳ LFU

این را هم مشابه MRU به سادگی از تغییر در LRU بدست می آوریم. در LRU ما هنگام قرار دادن در کش به آن یک امتیاز اضافه می کردیم اما اینجا امتیاز آن باید از صفر شروع شود و بعد از این در هر بار استفاده یکی افزوده شود. اینجا امتیاز ها نمایانگر تکرر استفاده هستند نه ترتیب آخرین استفاده.

```

1 #ifndef REPLACEMENT_LFU_H
2 #define REPLACEMENT_LFU_H
3
4 #include <vector>
5
6 #include "cache.h"
7 #include "modules.h"
8
9 class lfucache : public champsim::modules::replacement
10 {
11     long NUM_WAY;
12     std::vector<uint64_t> last_used_cycles;
13     uint64_t cycle = 0;
14
15 public:
16     explicit lfucache(CACHE* cache);
17     lfucache(CACHE* cache, long sets, long ways);
18
19     // void initialize_replacement();
20     long find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
21                     champsim::address full_addr, access_type type);
22     void replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
23                                access_type type);
24     void update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
25                                  access_type type, uint8_t hit);
26     // void replacement_final_stats()
27 };
28
29 #endif

```

فایل هدر

```

1 #include "lfu.h"
2
3 #include <algorithm>
4 #include <cassert>
5
6 lfucache::lfucache(CACHE* cache) : lfucache(cache, cache->NUM_SET, cache->NUM_WAY) {}
7
8 lfucache::lfucache(CACHE* cache, long sets, long ways) : replacement(cache, NUM_WAY(ways), last_used_cycles(static_cast<std::size_t>(sets * ways), 0)) {}
9
10 long lfucache::find_victim(uint32_t triggering_cpu, uint64_t instr_id, long set, const champsim::cache_block* current_set, champsim::address ip,
11                           champsim::address full_addr, access_type type)
12 {
13     auto begin = std::next(std::begin(last_used_cycles), set * NUM_WAY);
14     auto end = std::next(begin, NUM_WAY);
15
16     // Find the way whose last use cycle is most distant
17     auto victim = std::min_element(begin, end);
18     assert(begin <= victim);
19     assert(victim < end);
20     return std::distance(begin, victim);
21 }
22
23 void lfucache::replacement_cache_fill(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip, champsim::address victim_addr,
24                                       access_type type)
25 {
26     // Mark the way as being used on the current cycle
27     last_used_cycles.at((std::size_t)(set * NUM_WAY + way)) = 0;
28 }
29
30 void lfucache::update_replacement_state(uint32_t triggering_cpu, long set, long way, champsim::address full_addr, champsim::address ip,
31                                       champsim::address victim_addr, access_type type, uint8_t hit)
32 {
33     // Mark the way as being used on the current cycle
34     if (hit && access_type(type) != access_type::WRITE) // Skip this for writeback hits
35         last_used_cycles.at((std::size_t)(set * NUM_WAY + way))++;
36 }

```

فایل پیاده سازی

ARC ۴-۲-۳

سیاست جایگزینی ARC با توجه به مقاله اصلی این سیاست

Algorithm Cache Replacement Adaptive an with LRU Outperforming

نوشته شده است. در ادامه PseudoCode مربوطه آورده شده است:

ARC(c) INITIALIZE $T_1 = B_1 = T_2 = B_2 = 0, p = 0$. x - requested page.

Case I. $x \in T_1 \cup T_2$ (a hit in ARC(c) and DBL($2c$)): Move x to the top of T_2 .

Case II. $x \in B_1$ (a miss in ARC(c), a hit in DBL($2c$)): Adapt $p = \min\{c, p + \max\{|B_2|/|B_1|, 1\}\}$. REPLACE(p). Move x to the top of T_2 and place it in the cache.

Case III. $x \in B_2$ (a miss in ARC(c), a hit in DBL($2c$)): Adapt $p = \max\{0, p - \max\{|B_1|/|B_2|, 1\}\}$. REPLACE(p). Move x to the top of T_2 and place it in the cache.

Case IV. $x \in L_1 \cup L_2$ (a miss in DBL($2c$) and ARC(c)): case (i) $|L_1| = c$:
 If $|T_1| < c$ then delete the LRU page of B_1 and REPLACE(p).
 else delete LRU page of T_1 and remove it from the cache.
 case (ii) $|L_1| < c$ and $|L_1| + |L_2| \geq c$:
 if $|L_1| + |L_2| = 2c$ then delete the LRU page of B_2 .
 REPLACE(p).
 Put x at the top of T_1 and place it in the cache.

Subroutine REPLACE(p)
 if $(|T_1| \geq 1)$ and $((x \in B_2 \text{ and } |T_1| = p) \text{ or } (|T_1| > p))$ then move the LRU page of T_1 to the top of B_1 and remove it from the cache.
 else move the LRU page in T_2 to the top of B_2 and remove it from the cache.

این سیاست جایگزینی از دو لیست اصلی T_1 و T_2 و دو لیست "روح" B_1 و B_2 برای پیاده سازی الگوریتم خود استفاده می کند. اجتماع دو لیست اول حافظه نهان ما را تشکیل می دهند و مقادیر درایه های حافظه نهان را ضبط می کنند. دو لیست دوم نیز متناظر در ادامه حافظه اصلی می آیند و نقش آن ها تنها نگهداری metadata درایه های حافظه است (یعنی تنها آدرس و برچسب یک درایه). اگر قرار دهیم:

$$C = W \times S$$

where C is the cache size, W is the number of ways, and S is the number of sets.

آنگاه می توان قرارگیری این چهار لیست را بدین گونه تصور کرد که T_1 و T_2 در کنار همدیگر C ای به اندازه c تشکیل می دهند، و ابتدای B_1 به انتهای T_1 و ابتدای B_2 به انتهای T_2 متصل است. وظیفه T_1 نگه داشتن درایه هایی است که اخیراً یک بار استفاده شده اند. T_2 درایه هایی را نگه می دارد که اخیراً بیش از یک بار نگهداری شده اند. B_1 و B_2 هم به ترتیب برچسبی از درایه های خارج شده از T_1 و T_2 نگه می دارند. همچنین پارامتر p در این سیاست، هدف ما از نسبتی از C

که باید به T_1 مختص یابد را نگهداری می‌کند. در طول زمان و با پیشروی stream این نسبت بروز رسانی می‌شود تا جای بیشتری به لیست مهم تر اختصاص دهد.

در ادامه منطقی که با درخواست پردازنده از درایه x طی می‌شود را آورده ایم:

Algorithm 1 REPLACE(p)

if $|T_1| \geq 1$ and $((x \in B_2$ and $|T_1| = p)$ or $|T_1| > p)$

Move the LRU page of T_1 to the top of B_1

Remove it from the cache

else

Move the LRU page of T_2 to the top of B_2

Remove it from the cache

تابع بالا، منطق تخلیه کردن یک درایه در سیاست ARC است. هنگام پر بودن حافظه نهان اگر لیست T_1 نسبتی بیش از میزان دلخواه ما (p) از حافظه نهان را اشغال کرده است، پایین ترین عضو آن را حذف کرده و وارد لیست روح B_1 می‌کنیم. در غیر این صورت T_2 بیش از میزان مطلوب حافظه نهان را اشغال کرده است و باید پایین ترین درایه آن را حذف کرده و در لیست روح B_2 قرار دهیم. حال منطق اصلی سیاست:

۱. حالت نخست: x در T_1 یا T_2 قرار دارد

این به معنی یک hit است. درایه x باید از هر کجایی در این دو لیست که قرار دارد خارج شود و به ابتدای لیست T_2 برود، زیرا hit به این معنی است که این درایه اخیرا بیش از یک بار استفاده شده است.

۲. حالت دوم: x در B_1 قرار دارد این به معنی miss است. درایه x در حافظه نهان نیست، اما اخیرا در آن حضور داشته است. در این صورت باید نسبت مطلوب T_1 افزایش یابد - یعنی افزایش مقدار p . سپس تابع REPLACE صدا زده می‌شود و درایه x به آغاز لیست T_2 می‌رود

۳. حالت سوم: x در B_2 قرار دارد این به معنی miss است. به طور مشابه درایه x در حافظه نهان نیست، اما اخیرا در آن حضور داشته است. در این صورت باید نسبت مطلوب T_2 افزایش یابد - یعنی کاهش مقدار p . سپس تابع REPLACE صدا زده می‌شود و درایه x به آغاز لیست T_2 می‌رود

۴. حالت چهارم: x در هیچ یک از لیست ها قرار ندارد این به معنی miss است اما این بار درایه x

اخیرا استفاده نشده است. به کمک تابع *REPLACE* یک مکان در حافظه نهان خالی می‌شود و درایه x به دلیل این که تنها یک بار اخیرا استفاده شده است در اول لیست T_1 قرار می‌گیرد

بدین گونه سیاست ARC به شکل انعطاف پذیر ترکیبی از Least و Used Recently Least و Used Frequently را به ما ارائه می‌دهد که با پیشروی برنامه قدرت هر یک از این سیاست ها به کمک پارامتر p در حافظه نهان بیشتر یا کمتر می‌شود.

۵-۲-۳ LARC

برای این سیاست سعی کردیم مطابق شبه کد درون مقاله با داده ساختارهای صف گونه الگوریتم را پیاده سازی کنیم.

برای ذخیره داده های Q و Q_r برای تمام مجموعه ها، یک متغیرها را از جنس Vector تعریف می‌کنیم که خود حاوی چندین Vector است که هر کدام متناظر این صف ها برای یک مجموعه خاص هستند. از آنجایی که اعضای Q درون حافظه نهان وجود دارند، برای ذخیره سازی آن ها کافی است که اندیس آن ها درون مجموعه (یعنی مقدار way) را ذخیره کنیم. هر زمان که بلوک جدیدی وارد شد یا Hit اتفاق افتاد آن را به سر این صف می‌بریم. همچنین مقدار اولیه این لیست ها به دلیل پیاده سازی مستقیم تر حاوی تمامی اندیس ها می‌باشد.

برای Q_r نیز به شکل مشابه یک متغیر در کلاس تعریف می‌کنیم، با این فرق که اعضای آن باید نشان دهنده خود بلوک های حافظه باشند، برای این کار از داده ساختار `champsim::address` استفاده می‌کنیم. این نوع داده حاوی آدرس فیزیکی کامل است که به دلیل این که در حافظه نهان سطح آخر (LLC) هستیم بیت های آفست این آدرس ها صفر خواهند بود که یعنی هر آدرس نشان دهنده یک بلوک در حافظه اصلی خواهد بود.

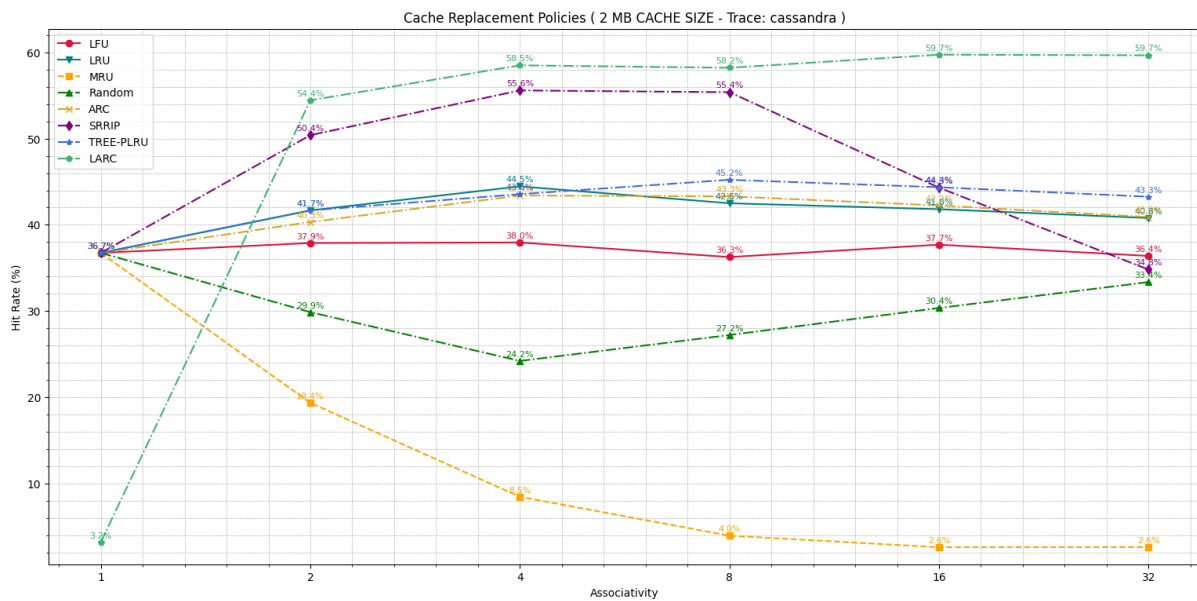
حال کافی است شروط الگوریتم را پیاده کنیم. شرط Hit در تابع `update_replacement_state` پیاده سازی شده است و Miss معادل تابع `find_victim` است. برای این که یک بلوک را اصلا وارد حافظه نهان نکنیم کافی است تعداد Way را به عنوان اندیس قربانی در این تابع خروجی بدهیم. همچنین نکته ای در اینجا وجود دارد آن هم این است که ممکن است نیاز به Write داشته باشیم اما این سیاست اجازه ورود به حافظه نهان را ندهد، در این صورت با خطا مواجه می‌شویم چرا که ساختار حافظه شبیه ساز ChampSim از نوع Write back است و ممکن نیست بدون ورود به حافظه نهان در بلوکی بنویسیم. برای حل این مشکل شرطی اضافه کردیم که در صورت Write همیشه انتقال به حافظه نهان اتفاق بیفتد.

۴ شبیه‌سازی و پارامترها

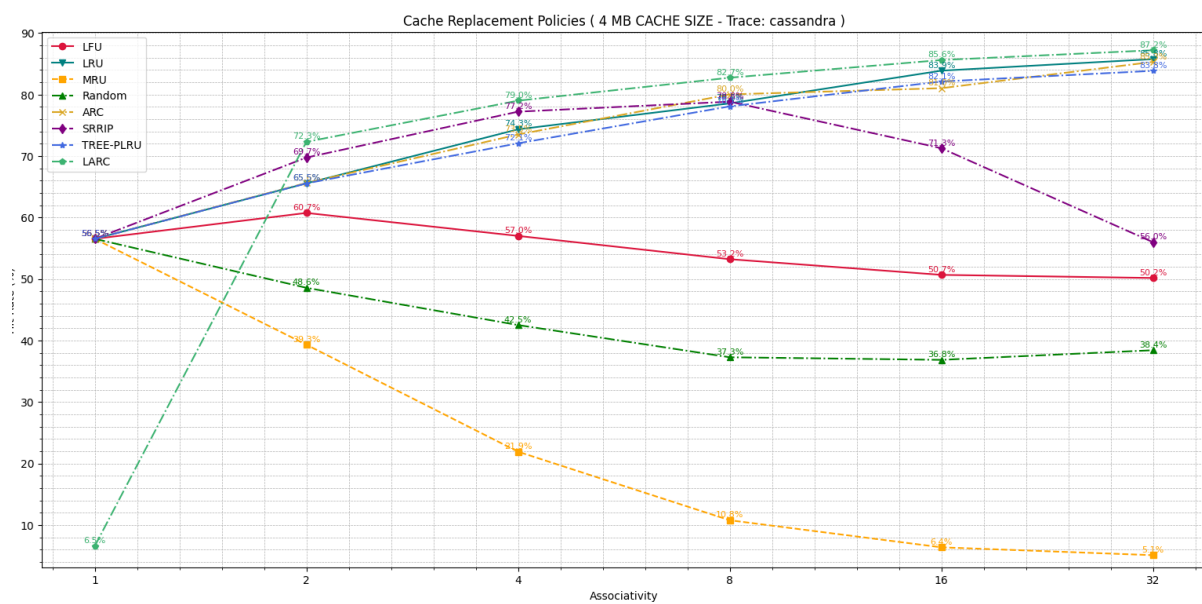
برای شبیه‌سازی و کسب داده موارد زیر در نظر گرفته شده است:

- Warmup Instructions: که تعداد دستورات گرم کردن است، در همه شبیه‌سازی‌های برابر 800,000 گرفته شده است.
- Simulation Instructions: که تعداد دستورات اصلی است، در تمام شبیه‌سازی‌های برابر 2,400,000 گرفته شده است.
- Associativity: که نشان‌دهنده تعداد بلوک در هر مجموعه است، برای بررسی عملکرد سیاست‌ها، این پارامتر را با مقادیر 1, 2, 4, 8, 16, 32 مقادیر کردیم تا نمودار نرخ برخورد را بر حسب این مقادیر رسم کنیم.
- Block Size: اندازه بلوک‌های حافظه نهان را در تمامی شبیه‌سازی‌ها برابر 64 بایت گرفتیم.
- Cache size(Set count): اندازه کل حافظه نهان لایه آخر معادل ضرب سایز بلوک در تعداد مجموعه و تعداد بلوک هر مجموعه است، شبیه‌سازی‌ها را در دو سری 2MB و 4MB برای اندازه حافظه اجرا کردیم که با تغییر تعداد مجموعه‌ها انجام شد. (در هر سری ضرب Associativity در تعداد مجموعه باید ثابت بماند)
- Trace: عبارت است از مجموعه دستوراتی که برای شبیه‌سازی استفاده می‌شوند. ما در این پروژه از دو مجموعه **Cassandra** و **Streaming** برای مقایسه مجموعه‌های مختلف استفاده شده است.

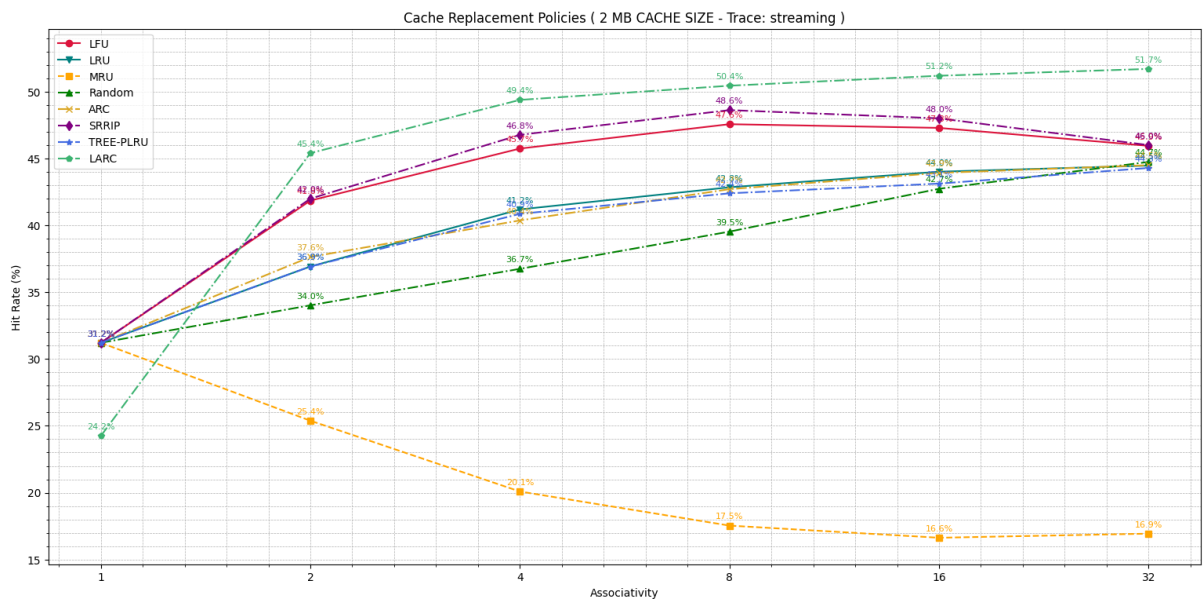
٥ تحلیل نتایج



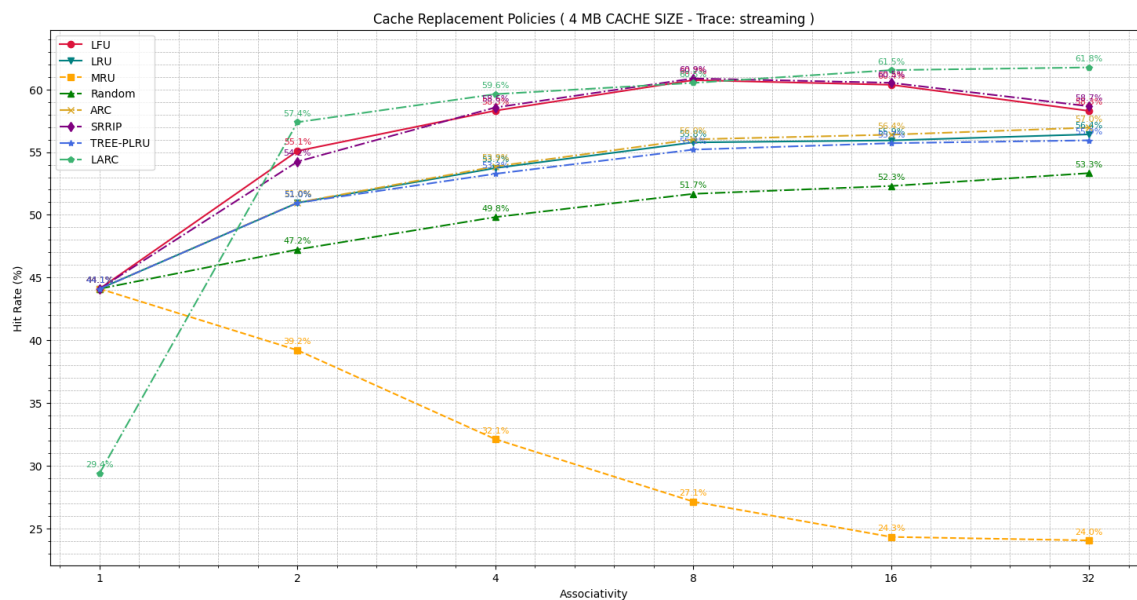
شکل ١: 2MB-cache size:cassandra



شکل ٢: 4MB-cache size:cassandra



شکل ۳: 2MB-cache size:streaming



شکل ۴: 4MB-cache size:streaming

ابتدا رابطه نرخ برخورد هرمنحنی با تعداد راه را جداگانه بررسی می کنیم و بعد به رابطه آن با بقیه می پردازیم. در نظر داشته باشید تحلیل دقیق ریاضی تحذب و تفرع و اثبات روابط از حوصله این مقاله خارج است و تاکید ما در این مقاله بر الگوریتم، پیاده سازی و شبیه سازی سیاست ها و ارزیابی آن ها براساس نرخ برخورد است. بنابراین در تحلیل نتایج نیز ما به دنبال قسمت ریاضیاتی اثبات مکشوفات خود نیستیم و صرفا به پیدا کردن روابط و دلیل شهودی آن ها بر اساس آموخته های خود می پردازیم.

توجه کنید way در ۱ چندان معنایی در جایگزینی ندارد برای همین عموما آن را نادیده می گیریم.

همچنین توجه کنید ما اختلافات نرخ برخورد کم و قابل چشم پوشی را تحت "هم قدرت" یا "هم رده" تحلیل می کنیم، این لغات معمولا وقتی به کار می روند که رتبه بندی ها دائما در تست های مختلف عوض می شود ولی اختلاف همانطور که گفته شده کم است.

ابتدا با MRU شروع می کنیم. این منحنی در تمامی حالات محدب بود و کمترین مقدار را دارد. این نشان می دهد این الگوریتم هرچند در تئوری در برخی حالات که درباره کد اجرایی اطلاعاتی داریم بهترین بازدهی را دارد اما در حالت کلی شکست خورده است.

سیاست Random بستگی به نوع دستورات می تواند محدب یا مقعر باشد پس نتیجه گیری خاصی نمی توان کرد، اما با این وجود رتبه دوم از آخر را دارد و همچنان از MRU بالاتر است. این نشان می دهد اگر هیچ فرضی راجع به دستورات نداشته باشیم همچنان از فرض بیشترین استفاده بازدهی بهتری دارد. همچنین در اکثر مواقع به بهترین نرخ برخورد بین همه سیاست ها از بدترین حالت نزدیک تر است. این مورد در کنار سادگی پیاده سازی و هزینه سخت افزاری می تواند گزینه مطلوبی برای طراحی ها و معماری های خاصی باشد.

سیاست بعدی که بررسی می کنیم SRRIP است. این سیاست همواره مقعر است و در نسخه ۸ مسیره پیشینه می شود. به طرز جالبی با LFU همگرا می شود، اما ممکن است تععرشان فرق کند که این نشان می دهد اگر تعداد مسیر ها کمی زیاد شود بازدهی این دو تفاوت کمی دارد. این سیاست تا قبل از way برابر ۸ رقیب جدی LARC است و از بقیه بالاتر است اما بعد از آن بازدهی چندانی ندارد و عملا در تعداد مسیر های محدودتر این روش برای تعیین امتیاز بهتر است. این نشان می دهد این نسخه که بر زمان استفاده تاکید دارد توانسته در مسیر های محدود از روش سنتی یعنی آخرین استفاده پیشین بگیرد. این توسعه مدیون مدیریت سن اعضا در طی زمان است.

حال که صحبت از LFU شد باید بگوییم هرچند همواره اکستریم گلوبال دارد اما لزوما محدب یا مقعر نیست و این به نوع دستورات بستگی دارد. به دلیل تفاوت های در دو تریس مورد استفاده می توانیم بفهمیم می تواند از LRU بهتر باشد اما این به شدت به نوع استفاده بستگی دارد.

LRU به طور کلی صعودی است و هم قدرت. ARC نسخه اصلاحی این الگوریتم یعنی Tree-PLRU تقریبا هم قدرت است و به نظر می رسد صرفا برتری در هزینه دارد. اما به این حقیقت توجه کنید که به طرز جالبی در way برابر ۲ دقیقا همان نرخ برخورد را دارد که این بخاطر این است که وقتی ۲ مسیر داریم درخت ما با یک بیت نشان داده می شود، و فلش جهت درخت همواره به سمت MRU می رود و نتیجتا همواره نتیجه یکسانی دارند در این حالت.

تمامی الگوریتم های پیشین معمولا از منطق مشابهی پیروی می کردند منتها یا داده ساختار متفاوتی استفاده می کردند یا نوع نگرششان به داده ساختار یکسان متفاوت بود. الان می خواهیم به الگوریتم ARC و نسخه اصلاحی آن یعنی LARC را بررسی کنیم. این الگوریتم ها تفاوت این الگوریتم

های پیشین را درک می کنند و تصمیم دارند با ایجاد یک فرایند پویا، در هر حالت بهینه ترین الگوریتم را انتخاب کنند و در صورت اشتباه خودشان را اصلاح کنند.

ARC معمولاً روند صعودی را دارد که یعنی با افزایش تعداد مسیر بازدهی اش بهتر می شود اما به طور کلی با LRU و Tree-PLRU هم قدرت است، به طور که اکثر مواقع بالاتر است. برعکس LRU که نسبتاً در یک ردیف قدرت با نسخه اصلاحی خود داشت، نسخه اصلاحی ARC به طرز چشمگیری از همه سیاست ها من جمله خود ARC بالاتر است و تابع صعودی دارد. این مورد را می توان در این تفاوت جستجو کرد که عمده تفاوت LRU و TPLRU در داده ساختار مورد استفاده است اما در ARC و LARC تفاوت مهم تری داریم. این تفاوت در مسئله مواجه با عدم برخورد است و نشان می دهد به طرز جالبی با هر نوع پارامتر گذاری در LARC همواره نتایج از همه سیاست ها بهتر است. این نشان می دهد که ایده ترکیب دو سیاست LRU و LFU پتانسیل زیادی برای بهبود دارد. نتایجی که ما به دست آوردیم با دانش تئوری ما و مقاله های مرتب مانند LARC و ARC همخوانی دارد، و دانش ما در عمل تایید شد. در تصویر زیر می توانید نتایج خود مقاله LARC را مشاهده کنید. همانطور که مشخص است نتایج ما نیز این بر نتایج این مقالات صحه می گذارد. طبیعتاً در مقادیر تفاوت داریم به دلیل تریس های متفاوت اما در همگی LARC برتری خود را ثابت کرده.

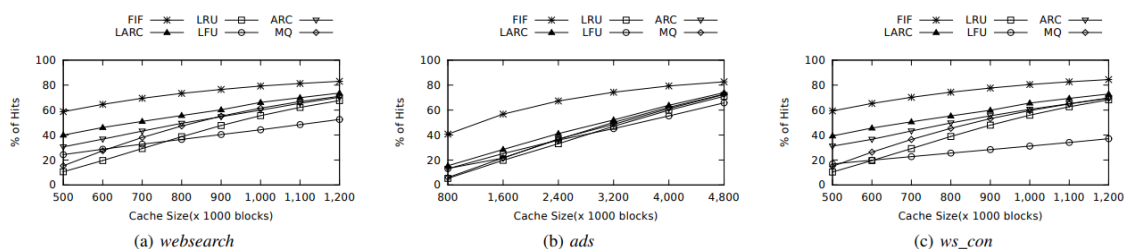


Fig. 4. Hit rate of different algorithms under read-dominant traces

حال با نتیجه گیری نهایی این قسمت را پایان می دهیم. شایان ذکر است که در این تحلیل و نتیجه گیری way برابر یک مورد توجه دقیقی قرار نگرفته، اما خوب است اشاره کنیم همواره بین همه نرخ یکسان است و دلیل یکسان نبودن LARC تفاوت در سیاست نوشتن است.

سیاست MRU بدترین سیاست در حالت کلی است. Random می تواند در شرایطی که بحث هزینه و سرعت مطرح است بهترین باشد. سیاست های LRU و Tree-PLRU تفاوت کمتری نسبت به ARC و LARC دارند. LARC بهترین سیاست در همه شرایط است. باقی سیاست ها همگرا می شوند، و معمولاً هم رده اند و در تریس های مختلف ممکن است عملکرد درخشان تری داشته باشند.

بهترین طراحی، ساده‌ترین راه‌حلی است که کار می‌کند.

— آلبرت اینشتین

در تمام این مقاله، هدف ما بررسی سیاست‌های مختلف و ارزیابی آن‌ها بر اساس نرخ برخورد بوده است. در این بخش می‌خواهیم مختصر به هزینه‌های پیاده‌سازی بپردازیم. از آنجا این موضوع هدف این مقاله نیست ما وارد جزئیات دقیق و ریاضیاتی آن‌ها نمی‌شویم بلکه به طور نسبی هزینه‌ها را بررسی می‌کنیم.

هزینه رندوم از همگی سیاست‌ها کمتر است. این به خاطر سادگی طراحی سخت افزاری است که عدد تصادفی تولید کند. پیچیدگی سخت افزار MRU و LRU از یک مرتبه است و به لحاظ هزینه معمولاً یک لیست برای نگهداری ترتیب نیاز داریم. مثلاً برای نگهداری ترتیب در یک مجموعه n بیتی می‌توانیم $\log n$ بیت برای ترتیب هر بلوک اختصاص دهیم که در می‌شود $n \log n$ اما در سیاست Tree-PLRU به دلیل استفاده از ساختار درختی هزینه نگهداری ما کاهش پیدا می‌کند. SRRIP نیز معمولاً با استفاده از شمارنده‌ها پیاده‌سازی می‌شود و از نظر اردر فضا در اردر LRU می‌شود (جدای هزینه خود الگوریتم در این درس نمی‌گنجد، اما به طور کلی باید از LRU بیشتر باشد). هزینه LFU نیز بستگی به طراحی دارد اما به طور کلی باز در اردر LRU است، این را در پیاده‌سازی ما هم می‌توانید ببینید که ساختار مشابهی دارند.

ARC و LARC به طور کلی هزینه پیاده‌سازی بیشتری دارند این موضوع در فایل‌های ما نیز مشهودند، خصوصاً که نیاز به ۴ صف دارند و باید اندیس آن‌هایی که حذف شده‌اند را هم به نوعی نگهدارند. پیچیدگی پیاده‌سازی این دو از همه بیشتر است و به سخت افزار بیشتری نیز نیاز دارد. به دلیل همین پیچیدگی می‌بایست به صورت نرم افزاری پیاده‌شوند و معمولاً برای هارد دیسک مورد استفاده قرار می‌گیرند.

References

- [1] Wikipedia contributors, “Adaptive replacement cache — wikipedia, the free encyclopedia,” 2024. Accessed: 2025-07-19.
- [2] N. Megiddo and D. S. Modha, “Arc: A self-tuning, low overhead replacement cache,” in *Proceedings of the 2nd USENIX Conference on File and Storage Technologies (FAST)*, USENIX Association, 2003.
- [3] S. Jiang and X. Zhang, “Making lru friendly to weak locality workloads: A novel replacement algorithm to improve buffer cache performance,” *IEEE Transactions on Computers*, vol.62, no.12, pp.2344–2357, 2013.
- [4] M. Hashemi, O. Mutlu, and B. Falsafi, “Accelerating dependent cache misses with an enhanced memory controller,” *IEEE Micro*, vol.36, no.4, pp.60–71, 2016.