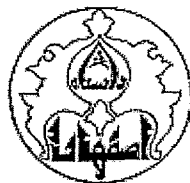


سید محمد



دانشگاه اصفهان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

پایان نامه‌ی کارشناسی ارشد رشته‌ی مهندسی کامپیوتر گرایش
معماری سیستم

طراحی یک سیستم حافظه نهان با مصرف انرژی بهینه با استفاده از تکنیک‌های
ویژه سخت افزاری و نرم افزاری

استاد راهنما

دکتر عباس وفایی

تأیید شده است
ششم مرداد

۱۳۸۸/۱۰/۲۷

پژوهشگر

نوید خیبر

خرداد ماه ۱۳۸۸

کلیه حقوق مادی مترتب بر نتایج مطالعات،
ابتکارات و نوآوری های ناشی از تحقیق
موضوع این پایان نامه متعلق به دانشگاه
اصفهان است.

شماره کارشناسی پایان نامه
رجاوت شده است
تحصیلات تکمیلی دانشگاه اصفهان



دانشگاه اصفهان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

پایان نامه ی کارشناسی ارشد رشته ی مهندسی کامپیوتر گرایش
معماری سیستم آقای نوید خیبر تحت عنوان

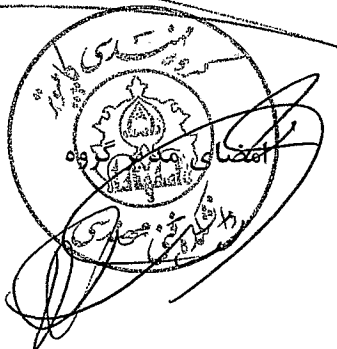
طراحی یک سیستم حافظه نهان با مصرف انرژی بهینه با استفاده از تکنیک های
ویژه سخت افزاری و نرم افزاری

در تاریخ ۱۳۸۸/۳/۳ توسط هیأت داوران زیر بررسی و با درجه عالی به تصویب نهایی رسید.

۱- استاد راهنمای پایان نامه دکتر عباس وفایی با مرتبه ی علمی استادیار امضا

۲- استاد داور داخل گروه دکتر محمدرضا خیام باشی با مرتبه ی علمی استادیار امضا

۳- استاد داور خارج از گروه دکتر سید مسعود سیدی با مرتبه ی علمی استادیار امضا



با سپاس فراوان از استاد گرامی و ارجمند
جناب آقای دکتر وفایی

تقدیم بہ

پرومادرم

چکیده

سرعت اجرای پردازنده های مدرن امروزی با سرعت چشمگیری در حال افزایش است که این امر موجب شده تا فاصله بین زمان دستیابی حافظه اصلی و سرعت اجرای پردازنده با نرخ بالاتری افزایش پیدا کند. یک راه حل مناسب جهت کاهش این فاصله استفاده از حافظه های نهان روی تراشه می باشد، بطوریکه هر روزه شاهد عرضه پردازنده هایی با مقادیر بیشتر حافظه نهان روی تراشه هستیم. اما با توجه به اینکه حدود ۴۰ درصد انرژی مصرفی پردازنده مربوط به حافظه نهان روی تراشه می باشد، افزایش میزان این حافظه موجب می شود تا توان مصرفی پردازنده به مقدار زیادی افزایش پیدا کند که این افزایش توان مصرفی پردازنده موجب می شود تا قابلیت اطمینان و چگالی تراشه نیز کاهش پیدا کند. علاوه بر این در سیستمهای قابل حمل نیز این افزایش توان مصرفی موجب کاهش عمر باتری می شود. اگرچه در پردازنده های مدرن امروزی اغلب از حافظه های نهان مجازی در سطح ۱ استفاده می شود، اما این نوع حافظه های نهان دارای یک مشکل عمده بنام مترادف می باشند. در این پایان نامه یک راهکار برای حل مشکل مترادف در حافظه های نهان مجازی و بهبود زمان دستیابی اصابت و توان مصرفی این نوع حافظه ها ارائه کرده ایم. سازمان حافظه نهان پیشنهادی ما، مجازی بوده و دارای مشکل مترادف نمی باشد. علاوه بر این، سازمان حافظه نهان پیشنهادی ما دارای زمان دستیابی پائین تر و توان مصرف کمتری نسبت حافظه های نهان نگاشت مستقیم و انجمنی- گروهی معادلش می باشد.

کلمات کلیدی: سیستم حافظه، حافظه نهان کم مصرف، حافظه نهان مجازی، مشکل مترادف

فهرست مطالب

صفحه	عنوان
	فصل اول مقدمه
۱-۱-۱	مقدمه..... ۱
۲-۱	حافظه نهان و جایگاه آن در سیستم کامپیوتری..... ۴
۳-۱	حافظه نهان یک منبع مناسب جهت کاهش مصرف توان پردازنده..... ۷
۴-۱	روشهای کاهش مصرف انرژی حافظه نهان..... ۸
	فصل دوم مفاهیم مقدماتی
۱-۲	انواع حافظه نهان از نظر سازمان..... ۱۰
۲-۲	ساختار داخلی حافظه نهان و روش دستیابی به آن..... ۱۲
۳-۲	انواع حافظه نهان از نظر نحوه دستیابی..... ۱۵
۴-۲	ساختار یک RAM ایستا و اجزای مختلف آن..... ۱۷
۱-۴-۲	عملیات نوشتن اطلاعات در SRAM..... ۱۸
۲-۴-۲	عملیات خواندن اطلاعات از SRAM..... ۱۹
۳-۴-۲	رمزگشای سطر..... ۲۴
۴-۴-۲	رمزگشای ستون..... ۲۴
۵-۲	مولفه های مختلف مصرف انرژی در حافظه نهان..... ۲۵
۶-۲	کارهای انجام شده جهت کاهش توان مصرفی حافظه نهان..... ۲۷
۱-۶-۲	بافر کردن بلوک..... ۲۷
۲-۶-۲	استفاده از زیر بانک در حافظه نهان..... ۲۹
۳-۶-۲	کاهش برچسب بصورت دینامیکی..... ۳۰
۴-۶-۲	آدرس دهی حافظه نهان با استفاده از کد گری..... ۳۱
۵-۶-۲	قطعه بندی Bit Line..... ۳۴
۶-۶-۲	پیش بینی راه..... ۳۵
۱-۶-۶-۲	حافظه نهان مرحله بندی شده..... ۳۶
۲-۶-۶-۲	حافظه نهان انجمنی- گروهی با قابلیت پیش بینی راه..... ۳۷
	فصل سوم مدل‌های تحلیلی تخمین زمان دستیابی و انرژی مصرفی
۱-۳	مدل تحلیلی تخمین زمان دستیابی..... ۳۹

۴۰	۱-۱-۳- سازمان حافظه نهان
۴۲	۲-۱-۳- مدل تحلیلی زمان دستیابی
۴۳	۱-۲-۱-۳- تاخیر رمزگشا
۴۴	۲-۲-۱-۳- تاخیر Word Line
۴۶	۳-۲-۱-۳- تاخیر بیت لاین/ حسگر- تقویت کننده
۵۰	۴-۲-۱-۳- تاخیر Data-Bus / Data-Out
۵۱	۵-۲-۱-۳- جمع بندی
۵۲	۲-۳- مدل تحلیلی تخمین انرژی مصرفی حافظه نهان
۵۲	۱-۲-۳- اتلاف انرژی در سلولهای SRAM
۵۳	۲-۲-۳- اتلاف انرژی در حافظه های نهان انجمی- گروهی
۵۵	۳-۲-۳- تعداد گذرها در حافظه های نهان انجمی- گروهی

فصل چهارم حافظه های نهان مجازی

۵۸	۱-۴- اهمیت و مزیت استفاده از حافظه های نهان مجازی
۶۰	۲-۴- مشکل همسانی در حافظه های نهان مجازی
۶۰	۱-۲-۴- مشکل مترادف
۶۱	۲-۲-۴- تغییرات نگاشت آدرس مجازی به فیزیکی
۶۲	۳-۴- همترازی آدرسهای مجازی و فیزیکی
۶۳	۴-۴- راهکارهای ارائه شده جهت حل مشکل مترادف
۶۴	۱-۴-۴- پیشگیری از مترادف
۶۴	۲-۴-۴- اجتناب از مترادف
۶۶	۳-۴-۴- تشخیص مترادف پویا - روش پایه
۶۷	۴-۴-۴- تشخیص مترادف پویا - نگاشتهای معکوس
۶۹	۵-۴- راهکارهای ارائه شده جهت حل مشکلات ناشی از تغییرات نگاشت صفحه
۶۹	۱-۵-۴- تغییرات نگاشت صفحه در حافظه های نهان V/P
۷۰	۲-۵-۴- تغییرات نگاشت صفحه در حافظه های نهان V/V
۷۰	۶-۴- مشکلات مختص به حافظه های نهان V/V
۷۰	۱-۶-۴- بیتهای حقوق دستیابی و بیتهای وضعیت صفحه

عنوان	صفحه
۴-۶-۲- دستیابی های خاص به TLB	۷۱
۴-۷-۷- بافر دم دستی مترادف (SLB)	۷۲
۴-۷-۱- آدرس مجازی اصلی بعنوان شناسه آدرس مجازی منحصر بفرد	۷۲
۴-۷-۲- طرح پایه : دستیابی موازی به حافظه نهان سطح ۱ و SLB	۷۳
فصل پنجم راهکار پیشنهادی	
۵-۱-۱- راهکار پیشنهادی: حافظه نهان چند قطعه ای	۷۶
۵-۲-۲- ارزیابی کارایی و توان مصرفی حافظه نهان چند قطعه ای با استفاده از مدل های تحلیلی	۸۵
۵-۳-۳- نتایج شبیه سازی	۸۸
۵-۴-۴- اضافه کردن قابلیت بافر کردن بلوک به حافظه نهان چند قطعه ای	۹۲
۵-۵-۵- نتایج شبیه سازی	۱۰۱
فصل ششم نتیجه گیری و کارهای آتی	
۶-۱-۱- نتیجه گیری	۱۰۴
۶-۲-۲- کارهای آتی	۱۰۵
پیوست الف آشنایی با نرم افزارهای شبیه ساز	
الف-۱- آشنایی با نرم افزار SimpleScalar	۱۰۷
الف-۱-۱- طریقه بدست آوردن و نصب ابزار	۱۰۸
الف-۲-۱- شبیه سازهای مجموعه ابزار	۱۱۰
الف-۲-۲- آشنایی با نرم افزار Wattch	۱۱۲
الف-۲-۱- روش مدلسازی	۱۱۲
الف-۳- آشنایی با نرم افزار HSPICE	۱۱۵
منابع و مآخذ	۱۱۷

فهرست شکل ها

عنوان	صفحه
شکل ۱-۱: افزایش چگالی توان مصرفی.....	۲
شکل ۲-۱: مثالی از یک سلسله مراتب حافظه چند سطحی.....	۴
شکل ۳-۱: مقایسه بین روند رشد سرعت دستیابی به حافظه با رشد کارایی پردازنده.....	۵
شکل ۱-۲: قالب آدرس برای دستیابی به حافظه نهان.....	۱۲
شکل ۲-۲: سازمان حافظه نهان نگاشت مستقیم.....	۱۳
شکل ۳-۲: سازمان حافظه نهان انجمنی- گروهی m-way.....	۱۴
شکل ۴-۲: ساختار داخلی یک حافظه RAM.....	۱۷
شکل ۵-۲: الف) سلول حافظه RAM ایستا ب) سلول حافظه RAM پویا.....	۱۸
شکل ۶-۲: سلول حافظه RAM ایستا به همراه مدارات پیش شارژ و حسگر/ تقویت کننده فلیپ فلاپ.....	۱۹
شکل ۷-۲: سلول حافظه RAM ایستا به همراه مدارات پیش شارژ و حسگر/ تقویت کننده تفاضلی.....	۲۲
شکل ۸-۲: حافظه RAM ایستا به همراه مدارات پیش شارژ و حسگر/ تقویت کننده.....	۲۲
شکل ۹-۲: دو نمونه ساده رمزگشای سطر.....	۲۴
شکل ۱۰-۲: یک نمونه ساده رمزگشای ستون.....	۲۵
شکل ۱۱-۲: تجزیه انرژی مصرف شده در حافظه نهان به مولفه های تشکیل دهنده آن.....	۲۷
شکل ۱۲-۲: تکنیک بافر کردن بلوک.....	۲۸
شکل ۱۳-۲: تکنیک تفکیک آرایه داده به تعدادی زیر بانک.....	۳۰
شکل ۱۴-۲: یک فرآیند نمونه و برچسب های مورد استفاده آن.....	۳۱
شکل ۱۵-۲: الف) ساختار اصلی یک ستون از سلولهای بیت ب) ساختار چند قطعه ای یک ستون از سلولهای بیت.....	۳۵
شکل ۱۶-۲: الف) حافظه نهان انجمنی- گروهی چهار راهه ب) حافظه نهان انجمنی- گروهی مرحله بندی شده چهار راهه.....	۳۶
شکل ۱۷-۲: الف) پیش بینی درست راه ب) پیش بینی نادرست راه.....	۳۷
شکل ۱-۳: الف) حافظه نهان انجمنی- گروهی دو راهه با $N_{dwl}=N_{dbl}=1$ ب) حافظه نهان انجمنی- گروهی دو راهه با $N_{dwl}=2, N_{dbl}=1$ ج) حافظه نگاشت مستقیم با $N_{dwl}=N_{dbl}=1$ د) حافظه نهان نگاشت مستقیم با $N_{dwl}=1, N_{dbl}=2$ (S بزرگ).....	۴۲
شکل ۲-۳: مراحل دستیابی به یک RAM.....	۴۳

- شکل ۳-۳: الف) مدل تاخیر Word Line قدیمی (ب) مدل تاخیر Word Line جدید..... ۴۵
- شکل ۳-۴: مدارات پیرامون سلول حافظه..... ۴۶
- شکل ۳-۵: الف) مدل تاخیر بیت لاین (ب) شکل موج بیت لاین..... ۴۷
- شکل ۳-۶: الف) مدل تاخیر بیت لاین (ب) شکل موج بیت لاین..... ۴۸
- شکل ۳-۷: مدارات مربوط حسگر/ تقویت کننده و درایور گذرگاه داده..... ۴۹
- شکل ۳-۸: مدارات مربوط درایور گذرگاه داده و درایور خروجی داده..... ۵۱
- شکل ۳-۹: الف) سلول SRAM شش ترانزیستوری (ب) مدارات پیش شارژ بیت لاین..... ۵۲
- ج) مدارات درایور Word Line..... ۵۲
- شکل ۴-۱: تغییرات نگاشت آدرس مجازی به فیزیکی..... ۶۲
- شکل ۴-۲: الف) ساختار حافظه نهان سطح ۱ و SLB..... ۷۴
- شکل ۵-۱: الف) قالب آدرس مجازی صادر شده توسط پردازنده (ب) قالب آدرس فیزیکی که موقعیت واقعی داده یا دستورالعمل را در حافظه مشخص می کند. (پ) قالب آدرسی که برای دستیابی به حافظه نهان از آن استفاده می شود..... ۷۷
- شکل ۵-۲: سازمان حافظه نهان چند قطعه ای برای یک کامپیوتر با ۶۴ کیلو بایت حافظه نهان داده سطح یک و ۶۴ کیلو بایت حافظه نهان دستورالعمل سطح یک و اندازه بلوک ۶۴ بایت و صفحات ۱۶ کیلو بایتی و گذرگاه آدرس ۶۴ بیتی الف) آرایه برجسب (ب) آرایه داده..... ۸۰
- شکل ۵-۳: مقایسه زمان دستیابی سمت داده..... ۹۰
- شکل ۵-۴: مقایسه زمان دستیابی سمت برجسب..... ۹۰
- شکل ۵-۵: انرژی مصرفی در Bit Line ها در اثر اجرای Compress..... ۹۰
- شکل ۵-۶: انرژی مصرفی در Bit Line ها در اثر اجرای Applu..... ۹۱
- شکل ۵-۷: انرژی مصرفی در Bit Line ها در اثر اجرای Wave5..... ۹۱
- شکل ۵-۸: انرژی مصرفی کل در اثر اجرای Compress..... ۹۱
- شکل ۵-۹: انرژی مصرفی کل در اثر اجرای Applu..... ۹۱
- شکل ۵-۱۰: انرژی مصرفی کل در اثر اجرای Wave5..... ۹۲
- شکل ۵-۱۱: حافظه نهان چند قطعه ای با قابلیت بافر کردن بلوک..... ۹۳
- شکل ۵-۱۲: سازمان حافظه نهان چند قطعه ای با قابلیت بافر کردن بلوک برای یک کامپیوتر با ۶۴ کیلو بایت حافظه نهان داده سطح یک و ۶۴ کیلو بایت حافظه نهان دستورالعمل سطح یک و اندازه بلوک ۶۴ بایت و صفحات ۱۶ کیلو بایتی و گذرگاه آدرس ۶۴ بیتی الف) آرایه داده (ب) آرایه برجسب..... ۹۸

۱۰۱.....	شکل ۵-۱۳ : انرژی مصرفی کل
۱۰۲.....	شکل ۵-۱۴ : مقایسه زمان دستیابی حافظه نهان داده
۱۰۲.....	شکل ۵-۱۵ : مقایسه زمان دستیابی حافظه نهان دستورالعمل
۱۰۶.....	شکل ۶-۱ : سازمان حافظه نهان ریز پردازنده AMD Opteron
۱۰۸.....	شکل الف-۱ : مرور گرافیکی مجموعه ابزار SimpleScalar
۱۱۳.....	شکل الف-۲ : ساختار کلی ابزار Wattch

فهرست جدول ها

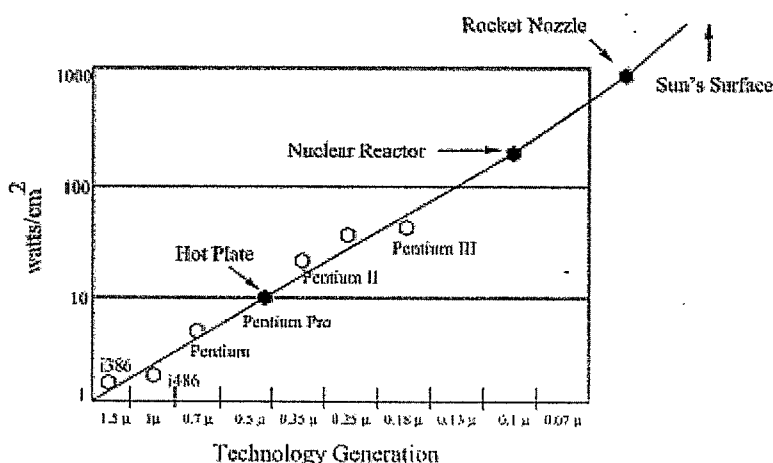
صفحه	عنوان
۳۲.....	جدول ۱-۲ : نمایش کد گری ۴ بیتی
۳۳.....	جدول ۲-۲ : مقایسه کلید زنی بیتی نمایش دودویی خالص و کد گری
۴۱.....	جدول ۳-۱ : پارامترهای استفاده شده در مدل تحلیلی تخمین زمان دستیابی
۷۲.....	جدول ۴-۱ : آدرسهای مجازی اصلی و ثانویه
۹۷.....	جدول ۵-۱ : حالات عملکرد حافظه نهان چند قطعه ای با قابلیت بافر کردن بلوک
۱۰۱.....	جدول ۵-۲ : پیکربندی های سیستم حافظه سناریوهای شبیه سازی
۱۱۱.....	جدول الف-۱ : پیکربندی پیش فرض حافظه نهان در sim-cache
۱۱۴.....	جدول الف-۲ : پارامترهای SPICE

فصل اول مقدمه

۱-۱- مقدمه

در گذشته وسایل الکترونیکی قابل حمل، شامل وسایلی مانند ساعت مچی و ماشین حساب بودند که امروزه این وسایل جای خود را به وسایل قابل حمل کاراتری مانند کامپیوترهای قابل حمل، گوشی های موبایل و دستگاه های پخش موزیک و پخش فیلم قابل حمل داده اند. از طرفی دیگر روز به روز شاهد تقاضاهای روز افزونی جهت ارائه خدمات چند رسانه ای بر روی این وسایل الکترونیکی قابل حمل هستیم، بنابراین این سیستمها باید دارای کارایی و سرعت پردازش بالایی باشند تا بتوانند خدمات مطلوب را به کاربر ارائه دهند. اما بالاتر رفتن سرعت پردازش، موجب می شود تا توان مصرفی نیز افزایش پیدا کند (بدلیل افزایش فرکانس سیگنال ساعت سیستم) و با توجه به اینکه در سیستمهای قابل حمل انرژی مورد نیاز سیستم بوسیله باتری مهیا می شود بنابراین توان مصرفی سیستم تاثیر مستقیم بر روی عمر باتری سیستم، و در نتیجه آن مدت زمانی که از سیستم می توان استفاده کرد، دارد. به همین دلیل عمر باتری، یکی از معیارهای مهم سنجش کیفیت چنین سیستمهایی بشمار می رود. از طرفی دیگر، توان مصرفی سیستم تاثیر مستقیم بر روی گرمای تولید شده توسط سیستم دارد بنابراین در سیستمهایی که توان مصرفی بالایی دارند باید تجهیزات خنک کننده زیادی جهت پائین نگه داشتن دمای سیستم به آن اضافه شود که این امر موجب افزایش حجم و هزینه ساخت سیستم می شود. علاوه بر این بدلیل مسائل

خنک کنندگی، قابلیت اطمینان و بسته بندی^۱، میزان توان مصرفی تاثیر تعیین کننده ای بر روی طراحی تراشه دارد و یک عامل محدود کننده برای تعداد ترانزیستوری است که می توان بر روی یک تراشه قرار داد. با پیشرفت تکنولوژی ساخت مدارات مجتمع و کوچکتر شدن اندازه ترانزیستورها، چگالی ترانزیستورهای روی تراشه افزایش پیدا کرده است و همانطوریکه در شکل ۱-۱ نشان داده شده هر چه چگالی ترانزیستورهای روی تراشه افزایش پیدا می کند، چگالی توان مصرفی و در نتیجه آن گرمای تولید شده توسط تراشه نیز افزایش پیدا خواهد کرد و در صورتیکه راه حل مناسبی برای کاهش توان مصرفی تراشه اندیشیده نشود، این عامل محدود کننده (بدلیل گرم شدن بیش از حد تراشه) باعث متوقف شدن روند پیشرفت تکنولوژی ساخت مدارات مجتمع می شود.



شکل ۱-۱: افزایش چگالی توان مصرفی [۱]

بنابراین با توجه به توضیحات فوق، جهت کاهش وزن و اندازه یک سیستم قابل حمل و افزایش عمر باتری آن، که از معیارهای مهم سنجش کارایی این نوع سیستمها می باشند، باید توان مصرفی این نوع سیستمها را کاهش دهیم. در سیستم های غیرقابل حمل نیز توان مصرفی معیار مهمی بشمار می آید زیرا هرچه توان مصرفی بیشتر باشد گرمای تولید شده توسط آن سیستم نیز بیشتر خواهد بود و سیستم به تجهیزات خنک کننده بیشتری نیاز خواهد داشت. بنابراین اگر در این سیستم ها توان مصرفی را پائین بیاوریم می توانیم بسیاری از تجهیزات خنک کننده را حذف کرده و قیمت سیستم را پائین بیاوریم و علاوه بر آن اندازه سیستم نیز می تواند کوچکتر شود.

^۱ Packaging

برای کاهش توان مصرفی یک سیستم لازم است تا مولفه های آن سیستم بطور دقیق بررسی شوند تا مشخص شود زمانیکه سیستم در حال فعالیت است هر مولفه آن چه مدت زمانی فعال است و هر مولفه در زمان فعالیتش چه مقدار انرژی در واحد زمان مصرف می کند، بدلیل اینکه مولفه های یک سیستم زمانی انرژی مصرف می کنند که فعال شده و مورد دستیابی قرار بگیرند این نوع انرژی مصرفی تحت عنوان انرژی مصرفی پویا^۱ شناخته می شود. علاوه بر این، نوعی دیگر از انرژی مصرفی بنام انرژی مصرفی ناشی^۲ نیز وجود دارد که این نوع انرژی بیانگر آن میزان انرژی می باشد که یک مولفه زمانیکه غیرفعال است مصرف می کند. امروزه تلاشهایی که برای کاهش انرژی مصرفی صورت می گیرد در جهت کاهش انرژی پویای یک سیستم می باشند [۲] که میزان این انرژی متناسب با تعداد باری است که مولفه های یک زیر سیستم مورد دستیابی قرار می گیرند. بنابراین اگر مدت زمان فعال بودن یا توان مصرفی دینامیکی مولفه های سیستم را کاهش دهیم، مصرف توان سیستم کاهش پیدا می کند. لذا در اولین گام لازم است تا بررسی کنیم کدامیک از مولفه های سیستم بیشتر مورد دستیابی قرار می گیرند و کدامیک دارای بیشترین توان مصرفی دینامیکی می باشند سپس در گام بعدی با کاهش توان مصرفی دینامیکی این مولفه ها یا کاهش تعداد باری که مورد دستیابی قرار می گیرند، می توان مصرف انرژی سیستم را تا حد زیادی کاهش داد.

یکی از مهمترین اجزاء مصرف کننده انرژی در یک سیستم کامپیوتری پردازنده می باشد بنابراین اگر بتوان پردازنده ای کم مصرف طراحی کرد می توان میزان توان مصرفی کل سیستم را تا حدود زیادی کاهش داد. تحقیقات زیادی بر روی طراحی پردازنده های کم مصرف انجام شده است و در حال حاضر نیز یکی از شاخه های فعال در تحقیقات می باشد. بدلیل اینکه اغلب دستیابی های پردازنده به حافظه نهان^۳ می باشد، برای کاهش توان مصرفی سیستم کامپیوتری استفاده از حافظه های نهان با مصرف انرژی پائین تبدیل به یک موضوع مهم در طراحی سیستمهای کامپیوتری مدرن شده است [۳].

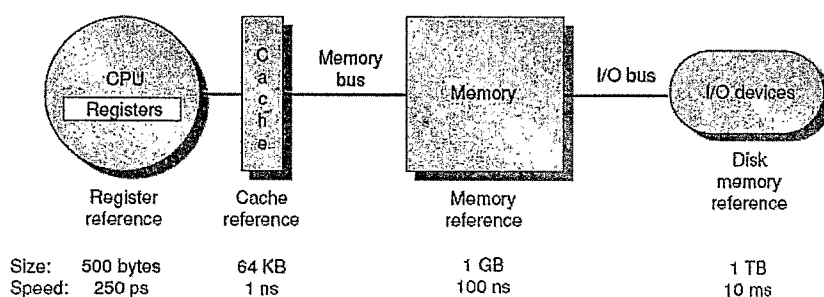
^۱ Dynamic Power

^۲ Leakage Power

^۳ Cache

۲-۱- حافظه نهان و جایگاه آن در سیستم کامپیوتری

پیشگامان علم کامپیوتر بدرستی پیش بینی کردند که برنامه نویسان در آینده، مقادیر زیادی از حافظه های سریع را طلب خواهند کرد. یک راه حل اقتصادی برای پاسخ به این خواسته استفاده از سلسله مراتب^۱ حافظه است، که از مزایای اصل محلی بودن ارجاعات^۲ و هزینه - کارایی^۳ تکنولوژیهای حافظه بهره برده و یک سیستم حافظه ارائه می دهد که سرعت آن تقریباً معادل سریعترین و گرانترین نوع حافظه و هزینه آن تقریباً معادل ارزانترین و کندترین نوع حافظه می باشد. اصل محلی بودن ارجاعات بیان می کند که اکثر برنامه ها کد و داده هایشان را به یک شکل مورد دستیابی قرار نمی دهند. محلی بودن در زمان^۴ و در فضای کد برنامه^۵ رخ می دهد. این مفهوم بعلاوه این نکته که سخت افزارهای کوچکتر سریعتر می باشند منجر به ایجاد یک سلسله مراتب حافظه براساس سرعت و اندازه متفاوت شد. در شکل ۲-۱ یک سلسله مراتب حافظه چند سطحی همراه با اندازه و سرعت دستیابی معمول برای هر سطح، نمایش داده شده است.



شکل ۲-۱: مثالی از یک سلسله مراتب حافظه چند سطحی [۵]

بدلیل اینکه حافظه های سریع گران هستند سلسله مراتب حافظه در سطوح مختلفی سازمان دهی شده است بطوریکه هر سطح آن، از سطح بعدی اش کوچکتر و سریعتر بوده و هزینه هر بایت آن گرانتر است.

بدلیل بهبود در کارایی پردازنده ها، سلسله مراتب حافظه اهمیت خاصی پیدا کرده است. شکل ۳-۱ مقایسه ای بین روند رشد کارایی پردازنده ها (سرعت اجرای پردازنده) و زمان دستیابی^۶ حافظه اصلی در زمان را

¹ Memory Hierarchy

² Locality

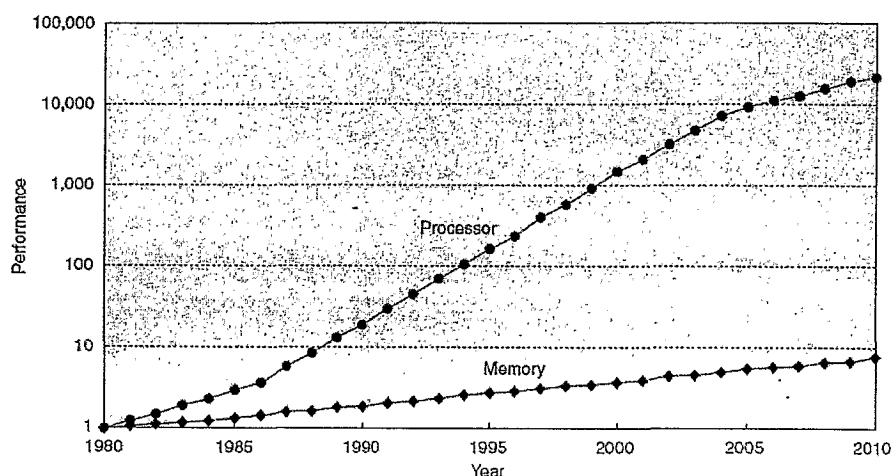
³ Cost-Performance

⁴ Temporal Locality

⁵ Spatial Locality

⁶ Access Time

نشان می دهد [۵]. بدلیل افزایش فرکانس سیگنال ساعت و استفاده از موازی سازی در سطح دستورالعمل^۱ (ILP)، پردازنده ها دارای نرخ بهبود کارایی ۶۰ درصد در سال هستند در حالیکه نرخ بهبود زمان دستیابی حافظه ها ۱۰ درصد در سال می باشد [۶].



شکل ۱-۳: مقایسه بین روند رشد سرعت دستیابی به حافظه با رشد کارایی پردازنده [۵]

واضح است که معماران کامپیوتر باید سعی کنند فاصله^۲ بین این دو متحنی را کاهش دهند. برای کاهش این فاصله سه راه حل وجود دارد:

۱. استفاده از حافظه های بزرگ، سریع و گرانقیمت
۲. استفاده از پردازنده هایی با کارایی در حد حافظه های کند و ارزان
۳. استفاده از سلسله مراتب حافظه و حافظه نهان سریع، کوچک و گرانقیمت

استفاده از روش اول باعث می شود تا سیستم، قیمت، مصرف انرژی و کارایی بالایی داشته باشد، استفاده از روش دوم باعث می شود تا سیستم، قیمت، مصرف انرژی و کارایی پائینی داشته باشد و روش سوم از مزایای هر دو روش بهره برده و معایب آنها را ندارد و باعث می شود تا سیستم قیمت و مصرف انرژی پائین، و کارایی بالایی داشته باشد.

حافظه نهان نامی است که به بالاترین (اولین) سطح از سلسله مراتب حافظه اختصاص داده می شود و آدرس صادر شده توسط پردازنده ابتدا با این سطح حافظه مواجه می شود. بدلیل اینکه استفاده از اصل محلی

¹ Instruction Level Parallelism

² Gap

بودن موجب بهبود کارایی می شود، از این اصل در سطوح مختلف بهره می برند بنابراین هر سیستم حافظه ای که داده های اخیراً استفاده شده را به قصد استفاده مجدد ذخیره می کند، تحت عنوان حافظه نهان شناخته می شود، بعنوان مثال حافظه های نهان فایل^۱ و حافظه های نهان نام^۲ نمونه هایی از این نوع حافظه ها می باشند.

زمانیکه پردازنده داده درخواست شده را در حافظه نهان پیدا می کند اصطلاحاً گفته می شود یک اصابت^۳ رخ داده است و داده موردنظر از حافظه نهان خوانده شده و در اختیار پردازنده قرار می گیرد. زمانیکه داده در حافظه نهان پیدا نشود اصطلاحاً گفته می شود یک عدم اصابت^۴ رخ داده است و در چنین حالتی، یک مجموعه داده با اندازه ثابت^۵ بنام بلوک^۶ که حاوی داده مورد نیاز نیز است، از حافظه اصلی خوانده شده و در حافظه نهان قرار می گیرد. بدلیل اینکه بنابر اصل محلی بودن در زمان، این احتمال وجود دارد که در آینده نزدیک، دوباره به این داده نیاز داشته باشیم و بنابر اصل محلی بودن در فضای کد برنامه، این احتمال نیز وجود دارد که در آینده نزدیک به داده های مجاور آن نیاز داشته باشیم، بنابراین در حالتی که عدم اصابت رخ میدهد، قرار دادن یک بلوک شامل داده و همسایه های آن از حافظه اصلی در حافظه نهان، برای بهبود کارایی پردازنده بسیار مفید می باشد. با توجه به اینکه اغلب دستیابی های پردازنده به حافظه نهان می باشد، بنابراین حافظه نهان نقش تعیین کننده ای در کارایی پردازنده دارد و بهبود کارایی حافظه نهان تاثیر مستقیم بر روی کارایی پردازنده دارد. به همین دلیل حافظه نهان از دیر باز مورد توجه معماران کامپیوتر بوده و سعی شده تا کارایی آن بهبود داده شود.

مدت زمان لازم جهت سرویس دهی به یک عدم اصابت، به تاخیر^۷ و پهنای باند^۸ حافظه وابسته است. تاخیر حافظه، بیانگر مدت زمان لازم برای بازیابی اولین کلمه از بلوک می باشد و پهنای باند حافظه، بیانگر مدت زمان لازم برای بازیابی بقیه بلوک را می باشد. عدم اصابت توسط سخت افزار سرویس دهی می شود و تا زمانیکه داده در دسترس قرار بگیرد، موجب توقف^۹ پردازنده هایی می شود که بصورت In-Order برنامه ها را اجرا می کنند. اما در پردازنده هایی که برنامه ها را بصورت Out-of-Order اجرا می کنند، عدم اصابت موجب متوقف شدن اجرای دستوراتی می شود که به داده درخواست شده نیاز دارند اما دستورات دیگر اجرا خواهند شد.

^۱ File Cache

^۲ Name Cache

^۳ Hit

^۴ Miss

^۵ Fixed-Size

^۶ Block

^۷ Latency

^۸ Band Width

^۹ Stall