

# بررسی تاثیر تنظیمات پارامترهای سخت‌افزاری بر انرژی مصرفی در الگوریتم ضرب برداری ماتریس‌های تنک بر روی پردازنده‌های گرافیکی

\* مینا عاشوری

\*\* فرشاد خون‌جوش

\* کارشناسی ارشد، مهندسی معماری کامپیوتر، دانشگاه شیراز، شیراز

\*\* استادیار، دانشکده مهندسی برق و کامپیوتر، دانشگاه شیراز، شیراز

تاریخ دریافت: ۱۳۹۷/۰۴/۰۷ تاریخ پذیرش: ۱۳۹۷/۰۹/۲۰

## چکیده

ضرب برداری ماتریس‌های تنک الگوریتمی ساده اما بخش بسیار مهمی از برنامه‌های جبر خطی و علمی در حوزه‌ی ریاضی و فیزیک است و به دلیل طبیعت قابل موازی سازی آن، پردازنده‌های گرافیکی یکی از گزینه‌های بسیار مناسب و مهم برای انتخاب بستر اجرایی آن است. در طی سال‌های اخیر با توجه به تاکید محققان برای در نظر گرفتن انرژی مصرفی به عنوان یکی از اهداف اصلی طراحی در کنار کارایی، تلاش‌های بسیار کمی جهت بهبود انرژی مصرفی این الگوریتم بر روی پردازنده‌ی گرافیکی انجام شده است. در این مقاله از منظر بهینگی مصرف انرژی در کارایی به دست آمده، به این مسئله پرداخته شده است.

با بهره‌وری از قابلیت تنظیم پیکربندی که در پردازنده‌های گرافیکی مدرن معرفی شده است، با بررسی آماری رفتار این الگوریتم هنگام استفاده از قالب‌های مختلف ذخیره سازی ماتریس تنک و تنظیمات مختلف سخت‌افزاری برای بیش از ۲۰۰ ماتریس نمونه‌ی تنک، بهترین تنظیمات پیکربندی برای الگوریتم ضرب برداری ماتریس تنک با قالب‌های مختلف ذخیره سازی بر روی پردازنده‌ی گرافیکی به دست آمده است. این پیکربندی برای هر قالب ذخیره سازی، به گونه‌ای انتخاب شده است که در تمام نمونه‌های بررسی شده به عنوان بهترین پیکربندی نتیجه داده باشد.

**واژه‌های کلیدی:** ضرب برداری ماتریس‌های تنک، انرژی مصرفی، کارایی، قالب‌های ذخیره سازی تنک، پردازنده‌ی گرافیکی.

## مقدمه

می‌شوند که شامل مقیاس بسیار بزرگی از برنامه‌های شبیه‌سازی مثل فیزیک انرژی بالا؛ شبیه‌سازی‌های علمی،

امروزه مسائل علمی و مهندسی مدرن، در طیف گسترده‌ای از برنامه‌های کاربردی روزانه و علمی با دقت بالا استفاده

<sup>1</sup> High Energy physics

همه منظوره<sup>۱</sup> است و از آن زمان تاکنون برای برنامه نویسی پردازنده‌های گرافیکی مورد استفاده‌ی همگانی قرار می‌گیرد. ضرب برداری ماتریس‌های تنک بخشی کوچک اما با اهمیت بسیار زیاد در زیربرنامه‌های اساسی تنک جبر خطی<sup>۱۲</sup> می‌باشد. بیشترین هزینه محاسباتی در بسیاری از روش‌های تکراری<sup>۳</sup> برای حل کردن سیستم‌های خطی در مقیاس بالا و مسائل با پاسخ خاص، مربوط به ضرب برداری ماتریس‌های تنک می‌باشد که در بسیاری از برنامه‌های کاربردی با اهمیت بالا مثل سیستم‌های شبیه سازی، تصویربرداری پزشکی، بازیابی اطلاعات، مدل‌سازی تغییرات اقلیمی و اقتصادی و بسیاری دیگر از طیف‌های برنامه‌های کاربردی وجود دارد. ماتریس‌های تنک متفاوتی که از برنامه‌های کاربردی متفاوت استخراج شده‌اند هرکدام دارای الگوی تنک بودن خاصی هستند و الگوی یکسانی ندارند. پس نمی‌توان از یک قالب یکسان برای ذخیره کردن همه استفاده کرد. ماتریس‌های متفاوت برای داشتن بهترین کارایی در عملیات ضرب برداری نیاز به قالب نمایش مناسب برای الگوی تنک بودن خود دارند. با توجه به جنبه‌ی بی‌قاعده‌ی این عملیات، هم در دستیابی به حافظه و هم در جریان کنترل<sup>۴</sup>، کارایی ضرب برداری ماتریس‌های تنک به هردو پارامتر الگوی تنک بودن ماتریس ورودی و ویژگی‌های سخت‌افزاری بستگی دارد. برای حل این مشکل تکنیک‌های بسیاری جهت بهینه سازی کارایی ضرب برداری ماتریس‌های تنک انجام گرفته است که شامل استفاده از قالب‌های متفاوت [۵][۶] و روش‌های متفاوت تنظیم خودکار<sup>۱۶</sup> با تنظیم کردن پارامترهای مختلف مطابق ساختار ماتریس است [۷]. البته اکثر تلاش‌ها برای بهبود ضرب برداری ماتریس‌های تنک در زمینه‌ی کارایی آن‌ها بوده و در حالی‌که این عملیات بر روی پردازنده‌های گرافیکی، با استفاده از قالب‌های

داده کاوی<sup>۲</sup> [۱]، جبرخطی [۲]، تجزیه و تحلیل گراف‌ها [۳]، پیش بینی آب و هوا و زلزله و ... می‌شوند که در همه‌ی آن‌ها، ماتریس‌های تنک کاربرد دارند. از دهه‌ی گذشته تا به امروز از روش محاسبات موازی استفاده می‌شود تا بتوان مسائل پیچیده را به روشی مقرون به صرفه حل کرد به طوری که بار کار<sup>۳</sup> بتواند بین پردازنده‌های مختلف تقسیم شود. در سال‌های اخیر، پردازنده‌های گرافیکی<sup>۴</sup> به دلیل هزینه‌ی پایین کارایی در واحد<sup>۵</sup> و میزان بالای کارایی در واحد توان<sup>۶</sup> به یکی از دستگاه‌های اصلی برای محاسبات با کارایی<sup>۷</sup> بالا تبدیل شده‌اند [۴]. قدرت محاسباتی بالا و پهنای باند این پردازنده‌ها موجب شده تا بتوان به عنوان جای‌گزینی برای سیستم‌های توزیع‌شده سنتی و یا پردازنده‌های چند هسته‌ای، در طیف وسیعی از شاخه‌های علوم مختلف از آن‌ها بهره جست. پردازنده‌های گرافیکی با استفاده از هزاران هسته<sup>۸</sup> و نخ<sup>۹</sup> حجم بسیاری از موازات را فراهم می‌آورند که سبب افزایش کارایی خواهند شد و برای برنامه‌های علمی موازی، سود کارایی با اجرا روی پردازنده‌های گرافیکی، به صورت مستقیم قابل دستیابی است. به همین دلیل، اخیراً محاسبات علمی، توسط جامعه‌ی پژوهش، برای اجرا به وسیله‌ی پردازنده‌های گرافیکی انطباق داده شده‌اند.

معماری یکپارچه‌ی دستگاه محاسبه<sup>۱۰</sup> (CUDA) به وسیله‌ی NVIDIA در سال ۲۰۰۶ به منظور بهره‌برداری از قدرت محاسباتی پردازنده‌های گرافیکی مطرح شد که یک زبان برنامه‌نویسی سطح بالا بر اساس پردازنده‌ی گرافیکی

1	General-purpose GPU	1
1	Basic sparse linear algebra (BLAS)	2
1	Iterative methods	3
1	Control flow	4
1	Format	5
1	Auto-tuning	6

2	Data mining
3	Work-load
4	Graphic processing units
5	Performance per unit
6	Performance per watt
7	High-performance computing
8	Core
9	Thread
10	Compute unified device architecture <sup>۱</sup>

هدف متضاد مصالحه‌ای ایجاد کند. در حالی که امروزه محققان تاکید می‌کنند بهینگی انرژی باید به عنوان یکی از اهداف اصلی در طراحی، در کنار کارایی قرار بگیرد، تلاش‌های بسیار کمی در جهت مطالعه‌ی مصرف انرژی در عملیات ضرب برداری ماتریس‌های تنک و بهبود آن انجام گرفته است.

هدف از انجام این تحقیق، بررسی تاثیر تنظیمات سخت‌افزاری بر روی انرژی مصرفی در واحد کارایی برای مسئله‌ی انتخاب قالب ذخیره سازی برای ضرب برداری ماتریس‌های تنک بر روی پردازنده‌های گرافیکی کپلر<sup>۲۲</sup> است. معیار بهینگی در این تحقیق مقدار کارایی در واحد توان مصرفی<sup>۲۳</sup> است که عبارت است از نسبت کارایی عملیات<sup>۲۴</sup> ضرب برداری ماتریس تنک به توان مصرف شده در این عملیات. به عبارتی دیگر عملیاتی با کارایی بالاتر و توان مصرفی کمتر مطلوب تر است. در این تحقیق مقدار کارایی معادل زمان خالص اجرای هسته‌ی ضرب برداری ماتریس تنک بر روی پردازنده‌ی گرافیکی تعریف شده و توان مصرفی هم بر حسب میلی وات<sup>۲۵</sup> توسط حسگرهای پردازنده‌ی گرافیکی کپلر اندازه گیری شده و بر حسب وات گزارش شده است.

در فصل ۲ به صورت خلاصه به پژوهش‌های پیشینی که در این زمینه انجام شده پرداخته می‌شود و در فصل ۳ پیش‌زمینه‌ای درباره‌ی پردازنده‌های گرافیکی، قالب‌های ذخیره‌سازی ماتریس تنک و زبان برنامه نویسی CUDA ارائه می‌شود. در فصل ۴ روش انجام این پژوهش توضیح داده شده و در فصل ۵ نتایج آزمایشات این پژوهش ذکر شده است. در فصل ۶ نیز نتیجه گیری و کارهایی که در آینده می‌توان با الهام از این پژوهش انجام شود بیان شده است.

ذخیره‌سازی مختلف، کارایی نسبتاً مناسبی دارد اما تلاش چندانی برای بهبود مصرف انرژی آن انجام نشده است. برای مثال در [8] برای تخمین کارایی و انتخاب بهترین قالب از مدل کردن کارایی استفاده شده است.

بهینگی در مصرف انرژی، در حالی که صنعت نیمه هادی از پردازنده‌های چند هسته ای<sup>۱۷</sup> به سمت پردازنده‌های با هسته‌های زیاد<sup>۱۸</sup> سوق پیدا می‌کند، به یک نیاز حیاتی در طراحی تبدیل شده است. طراحان در حال ترکیب ویژگی‌های سخت‌افزاری و نرم‌افزاری در این شتاب‌دهنده‌های بازده‌گرا<sup>۱۹</sup> هستند تا بتوان برنامه‌های کاربردی بدون ساختار با الگوهای متفاوت و بی‌قاعده‌ی دستیابی به حافظه را بر روی این پردازنده‌ها به خوبی اجرا کرد.

در حالی که مطالعات اخیر پیشنهاد می‌کند که بهینگی در مصرف انرژی به عنوان یک هدف اصلی در کنار کارایی در طراحی سخت‌افزار و نرم‌افزار قرار بگیرد، اغلب کارهایی که برای بهبود ضرب برداری ماتریس‌های تنک انجام گرفته در جهت بهبود کارایی به تنهایی بوده است و مصرف انرژی را شامل نشده است. مطالعات بسیار کمی جهت بررسی مصرف انرژی اجرای این هسته‌ی زوی پردازنده گرافیکی انجام شده است [9][10]. هیچکدام از این مطالعات، ارتباط بهینگی مصرف انرژی در این هسته را با الگوی تنک بودن ماتریس ورودی و ویژگی‌های سخت‌افزاری بررسی مطالعه نکرده‌اند. در [۱۱] به بررسی ارتباط ساختار ماتریس ورودی با انتخاب قالب‌های ذخیره‌سازی و تاثیر آن بر روی توان مصرفی پرداخته شده است ولی در نهایت تابه‌حال هیچ تلاشی انجام نشده است که دو هدف کارایی و بهینگی انرژی مصرفی را در طراحی یک سیستم بهینه برای ضرب برداری ماتریس‌های تنک بر روی پردازنده‌های گرافیکی دنبال کند و بین این دو

2 Trade-off	1
2 Kepler	2
2 Performance per watt	3
2 Operation	4
2 Milli-watt	5
2 Sensor	6
2 Watt	7

1 Multi-core processors	7
1 Many-core processors	8
1 Throughput-oriented accelerators	9
2 Kernel	0

سازی آن‌ها فقط از قالب CSR استفاده کرده‌اند. در [۱۴] انزت و همکاران مرزهای جدیدی را برای بهینگی انرژی و کارآیی برای محاسبات تنک بر روی ابرکامپیوترهای مبتنی بر پردازنده‌ی گرافیکی ارائه داده‌اند. LOBPCG به عنوان محک<sup>۳۰</sup> انتخاب شده است چون شامل عملیات ضرب برداری ماتریس‌های تنک نیز بوده است.

در باب اندازه‌گیری توان مصرفی در پردازنده‌های گرافیکی، نویسندگان [۱۵] مطالعه‌ای درباره‌ی چگونگی اندازه‌گیری توان لحظه‌ای و انرژی مصرفی با استفاده از حسگرهای روی برد پردازنده‌های گرافیکی جدید ارائه داده‌اند. برای پردازنده‌های گرافیکی بدون سنسور تعبیه شده می‌توان از مدل‌سازی‌هایی که برای توان صورت گرفته بهره برد یا از دستگاه‌های خارجی جهت اتصال به پردازنده‌ی گرافیکی و اندازه‌گیری توان مصرفی آن استفاده کرد [۱۶].

### ۳- مبانی نظری پژوهش

در این فصل به معرفی معماری پردازنده‌های گرافیکی و نحوه‌ی اجرای یک مدل برنامه‌نویسی رایج که همان CUDA است و همچنین عملیات ضرب برداری ماتریس-های تنک یا به اختصار<sup>۳۱</sup> SpMV و قالب‌های آن خواهیم پرداخت.

### ۳-۱ پردازنده‌های گرافیکی همه منظوره یا

#### GPGPU و مدل برنامه نویسی CUDA

مدل معماری GPGPU ارائه شده توسط شرکت NVIDIA مبتنی بر یک آرایه‌ای مقیاس پذیر از چندپردازنده‌هی جریان‌ی<sup>۳۲</sup> و چند نخ<sup>۳۳</sup> که هرکدام شامل تعدادی ثابت از پردازنده‌های عددی، یک یا چند واحد واکنشی، سخت افزار ویژه‌ی توابع خاص و حافظه‌ی سریع روی تراشه است که در نسل‌های کپلر و نسل‌های قبلی فرمی به حافظه‌ی نهان سطح یک<sup>۳۴</sup> و حافظه‌ی مشترک تقسیم می‌شود و در

### ۲- پژوهش‌های پیشین

یکی از اولین پیاده‌سازی‌های عملیات ضرب برداری ماتریس‌های تنک بر روی پردازنده‌های گرافیکی توسط بلتز و همکاران [۸] انجام گرفت. اما کار منحصر بفردی که در رابطه با سرعت بخشیدن به این هسته برای اجرا بر روی پردازنده‌های گرافیکی با قابلیت اجرای CUDA انجام گرفت توسط بل و گارلند در [۵][۱۲] انجام گرفت که در این پژوهش، آن‌ها یک مطالعه‌ی دقیق از قالب‌های ماتریس‌های تنک، الگوی دسترسی<sup>۳۵</sup> آن‌ها در پردازنده‌های گرافیکی و هسته‌ی عملیات‌های با قالب‌های کلاسیک و پیاده سازی شده با CUDA برای اجرا بر روی پردازنده‌های گرافیکی ارائه دادند. این قالب‌های کلاسیک عبارت بودند از COO, DIA, ELL, CSR و HYB.

تعدادی از پژوهش‌هایی که اخیرا انجام شده است با استفاده از قالب‌های زیادی که برای ماتریس تنک تعریف شده است سعی در بهبود کارآیی عملیات ضرب برداری ماتریس‌های تنک بر روی پردازنده‌ی گرافیکی داشته‌اند. اغلب آن‌ها به دنبال بهبود کارآیی بوسیله‌ی انتخاب قالب مناسب در زمان اجرای عملیات بوده‌اند. اما تلاش‌های بسیار کمی برای بهینگی مصرف انرژی هسته‌ی ضرب برداری ماتریس‌های تنک بر روی پردازنده‌های گرافیکی انجام شده است.

در [۱۳] نویسندگان میزان کارآیی در واحد وات را برای سه هسته از جمله هسته‌ی ضرب برداری ماتریس‌های تنک بر روی دو سکوی پردازنده‌ی Intel Sandy Bridge و پردازنده‌ی گرافیکی فرمی NVIDIA مطالعه کرده‌اند. آن‌ها نشان داده‌اند که با معیار کارآیی در واحد وات GFLOPs/watt این هسته بر روی پردازنده‌ی Intel Sandy Bridge بهتر عمل می‌کند ولی آن‌ها فقط از یک نوع ماتریس R-MAT استفاده کرده و برای ذخیره

<sup>3</sup> Benchmark 0

<sup>3</sup> Sparse matrix-vector multiplication<sup>1</sup>

<sup>3</sup> Streamed multi-processor 2

<sup>3</sup> Multi-threaded 3

<sup>3</sup> L1 cache 4

<sup>2</sup> Access pattern 8

<sup>2</sup> CPU 9

**قالب ELL**: ایده‌ی قالب ELL دسته بندی کردن همه‌ی اعداد غیر صفر ماتریس از سمت چپ آن و ذخیره‌ی ماتریس متراکم نتیجه است. این قالب از دو آرایه برای ذخیره‌ی المان‌های ماتریس متراکم و ذخیره‌ی شاخص ستونی این المان‌ها استفاده می‌کند.

**قالب SELL**: ایده‌ی آن بریدن ماتریس در تکه‌های هم طول ولی با عرض متفاوت است. این قالب از ۳ آرایه استفاده می‌کند. آرایه‌ی اول اشاره‌گرهای شروع هر Slice را ذخیره می‌کند. آرایه‌ی دیگر نیز مشابه قالب ELL است.

**قالب HYB**: برای ماتریسی که طول سطرهایشان متغیر است می‌توان از ترکیب دو قالب ELL و CSR برای ذخیره‌ی ماتریس استفاده کرد. زیرا برای سطرهایی با طول کمتر از طول ریسمان استفاده از قالب ELL باعث افزایش کارایی و برای سطرهایی با طول بزرگتر از ریسمان استفاده از قالب CSR باعث بهبود کارایی می‌شود. الگوریتم استفاده شده در این قالب تقسیم ماتریس به دو قسمت و ذخیره‌ی هر قسمت در هر کدام از این قالب‌ها است.

**قالب BCSR**: این قالب از سه آرایه تشکیل شده است. آرایه‌ی اول حاوی اشاره‌گرهایی به ابتدای هر ابر سطر است. با تقسیم ماتریس به بلوک‌های  $2 \times 2$  هر ۲ سطر در یک ابر سطر قرار می‌گیرند. آرایه‌ی دوم نیز مشابه CSR است اما در این قالب به جای ذخیره‌ی پشت سرهم المان‌ها در این آرایه، المان‌های هر بلوک به ترتیب بلوک‌ها بصورت پشت سرهم ذخیره می‌شوند. آرایه‌ی سوم نیز شاخص هر بلوک را در هر ابر سطر ذخیره می‌کند.

**قالب BELL**: این قالب یکی از نسخه‌های ELL است و در آن به جای ذخیره‌ی المان‌ها به صورت پشت سرهم، ماتریس به بلوک‌های هم اندازه تقسیم شده و بلوک‌های حاوی اعداد غیر صفر به صورت پشت سرهم مشابه قالب ELL ذخیره می‌شوند. هر بلوک و المان‌های درونش فقط به یک شاخص نیاز دارند پس دسترسی به حافظه کاهش می‌یابد. دسترسی به بردار ورودی برای همه‌ی المان‌های ابر سطر کاهش یافته و در بهترین حالت فقط یک بار اتفاق می‌افتد. این قالب از دو آرایه تشکیل شده است. آرایه‌ی اول شاخص ستونی هر بلوک را در ماتریس تقسیم شده به بلوک‌ها ذخیره می‌کند.

نسل‌های جدیدتر ماکسول و پاسکال به صورت فیزیکی جدا شده‌اند.

یک برنامه‌ی CUDA شامل یک برنامه‌ی میزبان و هسته است که برنامه‌ی میزبان روی CPU میزبان اجرا می‌شود و برنامه هسته روی دستگاه GPU اجرا می‌شود.

محاسبات توسط نخ‌ها انجام می‌شود که در بلوک‌ها گروه‌بندی می‌شوند. بیشتر از یک بلوک می‌تواند بر روی یک چند پردازنده‌ی یکسان اجرا شوند و بلوک‌های نخ‌ی به صورت هم‌رند اجرا می‌شوند.

طی فراخوانی هسته (شبکه‌<sup>۴۵</sup> نیز نامیده می‌شود)، برنامه‌ی میزبان تنظیمات اجرایی را که شامل موارد زیر است تعریف می‌کند:

- تعداد بلوک‌های نخ‌ی برای اجرا
- تعداد نخ‌ها در هر بلوک

## ۲-۳- قالب‌های ذخیره‌سازی تنک

در ضرب برداری ماتریس‌های تنک، هر المان ماتریس فقط یک بار مورد دسترسی قرار می‌گیرد. قالب‌های عمومی برای نگهداری و ذخیره کردن اطلاعات المان‌های غیر صفر ماتریس تنک محتمل سرباری زیادی می‌شوند. یکی از دلایل استفاده از قالب‌های تنک کاهش این سربار و کاهش دسترسی‌های غیر منظم، غیر تلفقی و کاهش واگرایی در نخ‌های یک ریسمان است. در پیاده‌سازی‌هایی که در این پژوهش انجام شده است برای ذخیره‌ی بردار ضرب شونده‌ی  $x$  از حافظه‌ی Texture بهره وری شده است.

**قالب CSR**: قالب CSR محبوب ترین قالب برای ذخیره کردن ماتریس تنک است. این قالب مقادیر غیر صفر و شاخص ستونی آن‌ها را به صورت غیر ضمنی در دو آرایه ذخیره می‌کند و از یک آرایه‌ی سوم برای ذخیره‌ی مرزهای هر سطر استفاده می‌کند. پیاده سازی این قالب را به دو صورت عددی و برداری می‌توان انجام داد که در نوع برداری مشهور به VCSR از روش parallel reduction استفاده می‌شود.

تعریف دقیق تعداد نخ‌ها در یک بلوک نقش اساسی در رسیدن به معیار کارایی در وات و اشغال<sup>۳۶</sup> پردازنده‌ی گرافیکی دارد. درصد اشغال که نرخ بهره‌وری از  $SM^2$  را نشان می‌دهد بسیار وابسته به تعداد نخ‌ها در بلوک نخ‌ی است زیرا برای تعداد بلوک فعال در هر  $SM$  محدودیتی وجود دارد و با تعریف مقدار صحیح اندازه‌ی این بلوک‌ها میتوان تعداد نخ‌های فعال در هر  $SM$  را به مقدار حداکثر آن نزدیک کرده و درصد اشغال  $SM$  و در نتیجه بهره‌وری از آن را افزایش داد. البته درصد اشغال به محدودیت منابعی دیگر نظیر تعداد ثبات‌ها برای هر نخ و حافظه‌ی مشترک موجود نیز بستگی دارد.

تعداد بلوک‌های فعال در هر  $SM$  به محدودیت منابعی مثل تعداد ثبات‌ها نیز بستگی دارد. در پردازنده‌های گرافیکی با قابلیت محاسباتی<sup>۳۸</sup> به بالا می‌توان به وسیله‌ی کامپایلر تعداد ثبات‌ها را از حداقل ۱۶ تا حداقل ۲۵۵ تعریف کرد. کاهش تعداد ثبات‌ها برای هر نخ تا تعداد مورد نیازشان می‌تواند باعث شود که محدودیت ثبات‌های موجود کمتر شود و نخ‌های فعال بیشتری برای اجرا روی پردازنده‌ی گرافیکی وجود داشته باشد و موازی بودن در سطح نخ<sup>۳۹</sup> بیشتری ایجاد کرد. مقدار ثبات‌ها برای هر نخ را می‌توان با پرچم `maxrregcount` در زمان کامپایل تعریف کرد.

قابلیت تنظیم سلسله مراتب حافظه در پردازنده‌های گرافیکی باعث شده است که برنامه‌های کاربردی محدود به حافظه نتوانند کارایی و توان مصرفی بهتری بر روی پردازنده‌ی گرافیکی داشته باشند. پردازنده‌های گرافیکی نسل‌های قبل فقط دارای حافظه‌های ثابت<sup>۴۰</sup> و بافته بودند اما پردازنده‌های گرافیکی نسل‌های جدید دارای قابلیت‌های نهان سازی سطح اول و دوم یا همان  $L1$  cache و  $L2$  chache نیز هستند. اندازه‌ی حافظه‌ی نهان سطح اول یا

آرایه‌ی دوم المن‌های بلوک‌ها را پشت سرهم و مشابه قالب  $ELL$  ذخیره می‌کند.

قالب  $DIA$  همان طور که از نام قالب  $DIA$  که مخفف قطری است مشخص است، این قالب مناسب ماتریس‌هایی با ساختار قطری است زیرا قطرهای ماتریس را به صورت ستونی در یک ماتریس فشرده ذخیره می‌کند.

قالب  $CDS$ . این قالب جهت بهره‌وری از ساختار بلوکی در ماتریس‌های محاسبات شبکه‌ای طراحی گردیده است. با تقسیم ماتریس به بلوک‌های هم اندازه، در این قالب به جای ذخیره‌ی هر قطر غیر صفر ماتریس، هر ابر قطر ذخیره می‌شود.

#### ۴- روش‌شناسی

در [۱۷] نشان داده شده است که انتخاب قالب مناسب برای ذخیره‌سازی ماتریس تنک تاثیر بسیار زیادی بر روی کارایی خواهد داشت. از طرفی مصرف انرژی در پردازنده‌ی گرافیکی در رابطه‌ی مستقیم با کارایی است پس می‌توان گفت که انتخاب قالب متناسب با ساختار ماتریس ورودی، می‌تواند با بهبود دسترسی به حافظه، کاهش سربار محاسباتی و استفاده از پهنای باند حافظه، استفاده‌ی بهینه مناسب از حافظه‌ی نهان و ... باعث بهبود در مصرف انرژی شود. علاوه بر انتخاب قالب مناسب، پارامترهای قابل تنظیمی در پردازنده‌ی گرافیکی وجود دارند که نتایج این پژوهش و پژوهش‌های دیگر نشان داده‌اند که با تنظیم صحیح آن‌ها می‌توان تاثیر مثبتی در بهبود کارایی و توان مصرفی ایجاد کرد. این پارامترهای سخت افزاری عبارتند از:

۱. اندازه‌ی بلوک نخ‌ی که تعداد نخ‌های هر بلوک را مشخص می‌کند

۲. تعداد ثبات‌های اختصاص داده شده به هر نخ

۳. سلسله مراتب حافظه

اندازه‌ی بلوک نخ‌ی پارامتری است که به شدت بر روی کارایی و توان مصرفی پردازنده‌ی گرافیکی تاثیر می‌گذارد.

<sup>3</sup> Occupancy

6

<sup>3</sup> Streaming multiprocessor 7

<sup>3</sup> Compute capability 8

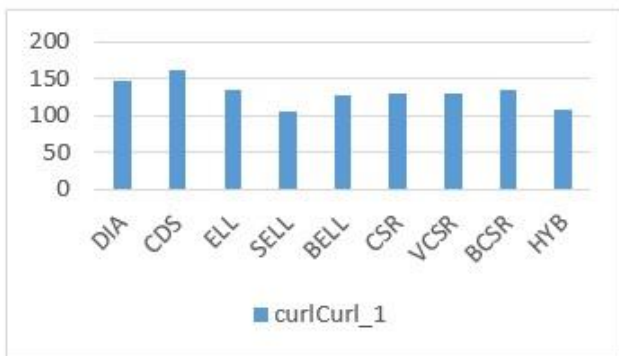
<sup>3</sup> Thread-level parallelism 9

<sup>4</sup> Flag 0

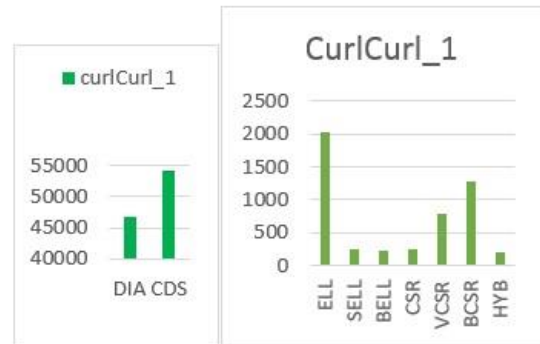
<sup>4</sup> Constant 1

سخت‌افزاری مناسب را برای پردازنده‌ی گرافیکی انتخاب کرد. این ایده بر اساس آزمایشاتی شکل گرفته است که نشان داده‌اند برای هر ماتریس تنک ورودی با هر الگوی ساختاری تنظیمات سخت‌افزاری و قالب خاصی برای ذخیره‌سازی وجود دارد که بهترین میزان کارایی در وات را نتیجه می‌دهد.

پس در ابتدا نیاز به انجام آزمایشاتی است که نحوه‌ی تاثیر تنظیم این پارامترها پس از انتخاب قالب مناسب را بر کارایی و توان مصرفی نشان دهد. برای انجام این آزمایشات تعدادی ماتریس تنک از مجموعه‌ی ۲۰۰ ماتریس تنک ورودی که از کتابخانه‌ی تنک دانشگاه فلوریدا جمع‌آوری شده و به صورت تصادفی انتخاب شده به طوری که اندازه‌ی ماتریس متراکم شده‌ی آن‌ها در قالب‌های این پژوهش، در حافظه‌ی RAM سیستم قابل تعریف باشد. محور افقی نمودارها نشان‌دهنده‌ی مقدار کارایی در وات با تنظیمات سخت‌افزاری است و  $t/p$  و  $perf/pow$  مخفف  $Performance(time)/power$  می‌باشد و اعداد بعد به ترتیب اندازه‌ی بلوک، تعداد ثابت در نخ و اولویت سلسله مراتب حافظه را نشان می‌دهد.



L1 cache و حافظه‌ی اشتراکی توسط برنامه نویس قابل تنظیم است. عملیات ضرب برداری ماتریس تنک نیز یک برنامه‌ی محدود به حافظه است و معمولاً به دلیل دسترسی-هایس به حافظه‌ی سراسری، تاخیر اجرایی و متوسط توان مصرفی نیز به تبع آن افزایش پیدا می‌کند. با معرفی پردازنده‌های گرافیکی مبتنی بر حافظه‌ی نهان و قابلیت تنظیم پیکربندی حافظه‌ی اشتراکی و نهان، می‌توان این نوع دسترسی را مدیریت کرد. برای برنامه‌های با الگوی دسترسی منظم به حافظه بهتر است اندازه‌ی حافظه‌ی اشتراکی را افزایش داد و برای برنامه‌های با الگوی دسترسی نامنظم بهتر است اندازه‌ی حافظه‌ی نهان را بیشتر کرد تا بهره‌وری از این حافظه با توجه به الگوی دسترسی نامنظم افزایش یابد [۱۸]. هدف این پژوهش یافتن تنظیمات سخت‌افزاری مناسب برای قالب‌های ذخیره‌سازی تنک در الگوریتم ضرب برداری ماتریس تنک، با انرژی مصرفی بهینه در واحد کارایی است. معیار بهینگی در این سیستم برای یک عملیات ضرب برداری معیار کارایی در واحد توان یا کارایی در وات است.



شکل ۱. زمان اجرای الگوریتم ضرب برداری ماتریس تنک CurlCurl\_1 برحسب میکروثانیه در قالب‌های DIA, CDS, ELL, SELL, BELL, CSR, VCSR, BCSR, HYB

شکل ۲. توان مصرفی الگوریتم ضرب برداری ماتریس تنک CurlCurl\_1 برحسب وات در قالب‌های DIA, CDS, ELL, SELL, BELL, CSR, VCSR, BCSR, HYB

۱-۴- معیار کارایی در واحد توان مصرفی (وات)

در این آزمایشات ۹ قالب ماتریس تنک مورد بررسی قرار گرفته‌اند. نتیجه‌ی آزمایشات در نمودارهای زیر نشان می‌دهند که متناسب با ساختار ماتریس ورودی، انتخاب یک

به این معنی که تنظیماتی مطلوب‌تر است که کمترین میزان مصرف انرژی را به ازای کارایی بیشتر نتیجه دهد. برای دستیابی به بالاترین میزان این معیار برای هر ماتریس تنک ورودی، می‌توان متناسب با الگوی ساختاری ماتریس و قالب منطبق بر این ساختار برای فشرده‌سازی تنظیمات مناسب

برای اندازه‌گیری زمان اجرای الگوریتم از روش چند نخ‌ی استفاده شده است. در نخ والد و اصلی، یک حلقه از شروع ارسال هسته بر روی پردازنده‌ی گرافیکی تا اتمام انجام عملیات هسته مشغول نمونه‌گیری از توان مصرفی لحظه‌ای پردازنده‌ی گرافیکی و محاسبه‌ی مدت زمان اجرا است. برای نمونه‌گیری از حسگرهای پردازنده‌ی گرافیکی برای اندازه‌گیری توان مصرفی لحظه‌ای، از رابط کتابخانه‌ی<sup>۴۲</sup> NVML استفاده شده است. به دلیل محدودیت حسگرهای اندازه‌گیری توان مصرفی که روی پردازنده‌ی گرافیکی قرار دارند هر الگوریتم ۵۰۰ بار اجرا شده است و زمان اجرا و توان مصرفی در این پنجره اندازه‌گیری شده است. در این پنجره هرگاه توان مصرفی پایدار و به صورت دوره‌ای تکرار شد، متوسط توان لحظه‌ای‌هایی که در این دوره وجود دارند به عنوان متوسط توان مصرفی لحظه‌ای الگوریتم در نظر گرفته می‌شود.

#### ۲-۴- تنظیمات سیستم انجام آزمایشات

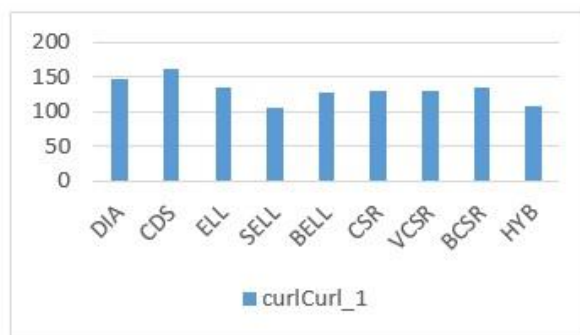
این پژوهش بر روی سیستمی انجام شده است که دو پردازنده‌ی K40m داشته است. تعداد نخ‌ها در هر SMX ۲۰۴۸ عدد بوده و دارای ۶۴ کیلوبایت حافظه‌ی روی تراشه است که به به ۳ صورت ۱۶ و ۴۸، ۳۲ و ۳۲ یا ۴۸ و ۱۶ کیلوبایت برای حافظه‌ی نهان و اشتراکی قابل تعریف است. پهنای باند این پردازنده‌ی گرافیکی به حافظه‌ی سراسری ۲۸۸ گیگابایت در ثانیه می‌باشد.

#### ۵- آزمایشات و نتایج

سه پارامتر سخت افزاری قابل تنظیم منتخب در این پژوهش، اندازه‌ی بلوک نخ‌ی، تعداد ثبات‌های اختصاص داده شده به هر نخ و سلسله مراتب حافظه است. تنظیمات قابل اعمال برای اندازه‌ی بلوک نخ‌ی از ۱۲۸ تا ۱۰۲۴ است و برای تعداد ثبات‌های اختصاص داده شده برای هر نخ مقادیر ۱۶، ۳۲، ۶۴ و ۲۵۶ در نظر گرفته شده و سلسله مراتب حافظه مطابق قابلیت تنظیم پیکربندی حافظه‌ی روی تراشه اعمال شده است. اندازه‌های ۱۶-۴۸ یا ۳۲-۳۲ کیلوبایت

قالب صحیح برای فشرده سازی و ذخیره‌ی ماتریس، تاثیر بسزایی بر روی کارایی و توان مصرفی و متناظر با آن معیار کارایی در واحد توان یا کارایی در وات خواهد داشت. با انتخاب قالب مناسب می‌توان با کاهش صفرهای افزوده متناظر با الگوی ساختاری ماتریس و کاهش فضای مورد نیاز در حافظه برای ذخیره‌ی ماتریس از پهنای باند حافظه بهره‌وری بهتری داشت. با انتخاب صحیح روش فشرده سازی همچنین می‌توان از محلیت حافظه‌ی نهان بهتر استفاده کرد و با ایجاد تعادل بار در نخ‌ها، از قدرت محاسباتی پردازنده‌ی گرافیکی استفاده‌ی بهینه‌تری داشت.

شکل ۱ نشان دهنده‌ی زمان اجرای هسته‌ی این الگوریتم با ماتریس تنک CurlCurl\_1 بر حسب میکرو ثانیه در این ۹ قالب است. شکل ۲ میزان توان مصرفی را برای این ماتریس در قالب‌های مختلف نشان می‌دهد. بهترین میزان کارایی برای این ماتریس در قالب‌های BELL، SELL و HYB به دست می‌آید ولی بهترین میزان توان مصرفی برای این ماتریس در قالب SELL است. اما شکل ۳ نشان می‌دهد که بهترین میزان کارایی در واحد توان مصرفی برای این ماتریس در قالب HYB است. لذا معیار کارایی یا انرژی مصرفی به تنهایی نمی‌تواند مناسب بودن یک قالب را نشان دهد و معیار مناسب میزان انرژی مصرف شده برای رسیدن به کارایی به دست آمده است.



شکل ۳. مقدار کارایی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس تنک CurlCurl\_1 در قالب‌های DIA, CDS, ELL, SELL, BELL, CSR, VCSR, BCSR, BCSR, HYB

برای اندازه‌گیری کارایی از تعریف معکوس زمان اجرا الگوریتم بر روی پردازنده‌ی گرافیکی استفاده شده است و

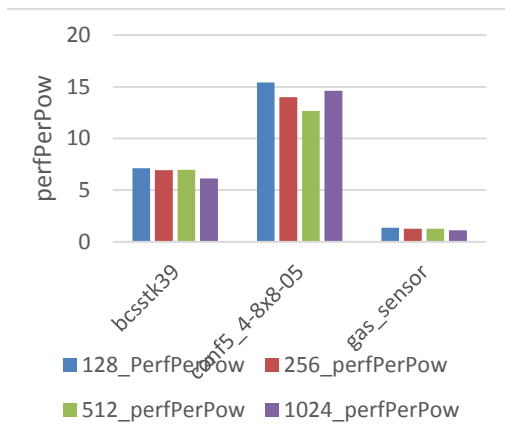
<sup>4</sup> Nvidia management library

برای حافظه‌ی اشتراکی و حافظه‌ی نهان قابل تخصیص است.

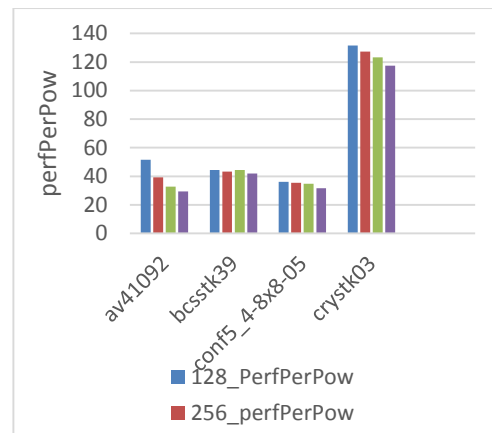
شکل‌های ۴ تا ۶ نشان‌دهنده‌ی تاثیر تنظیمات مختلف اندازه‌ی بلوک نخی بر روی کارایی در واحد توان برای چند ماتریس نمونه‌ی تنک در قالب‌های VCSR، DIA و BCSR است.

شکل‌های ۷ تا ۹ نشان‌دهنده‌ی تاثیر تنظیمات مختلف تعداد ثبات در نخ بر روی کارایی در واحد توان برای قالب‌های BELL، CDS و ELL می‌باشد. در محور افقی عدد اول، اندازه‌ی بلوک نخی و عدد دوم، تعداد ثبات‌های

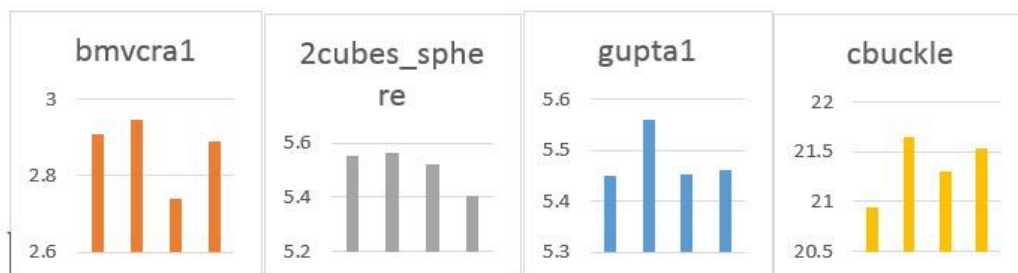
اختصاص داده شده به هر نخ را نشان می‌دهد. شکل ۱۰ نیز نشان‌دهنده‌ی تاثیر تنظیمات مختلف سلسله مراتب حافظه بر روی کارایی در واحد توان برای قالب VCSR است. نموداری که با حروف 11 پایان میابد مقدار کارایی در وات را برای تنظیماتی با اندازه‌ی حافظه نهان ۴۸ کیلوبایت نشان می‌دهد و نمودارهای دارای حروف eq و smem نیز نشان‌دهنده‌ی کارایی در وات برای تنظیمات اندازه‌ی حافظه‌ی نهان ۳۲ و ۱۶ است.



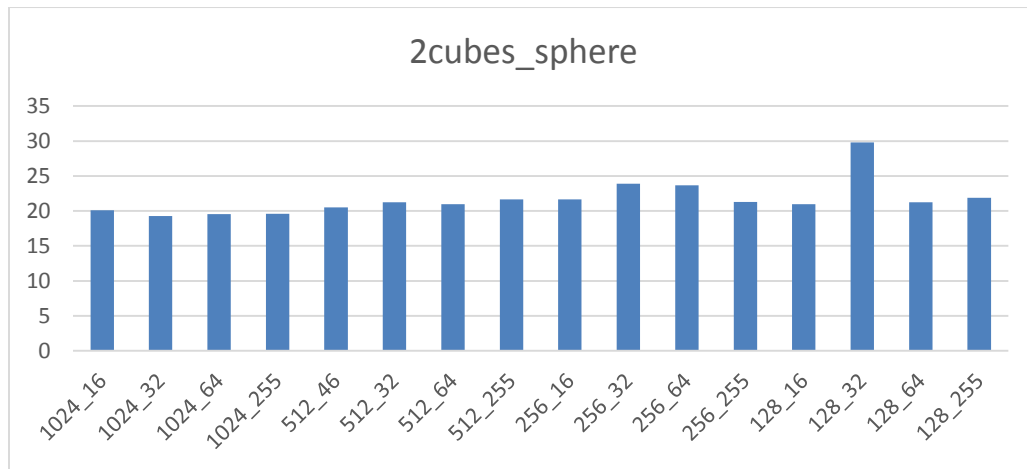
شکل ۵. مقدار کارایی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس تنک نمونه در قالب DIA در اندازه‌های بلوک نخی ۱۲۸، ۲۵۶، ۵۱۲ و ۱۰۲۴



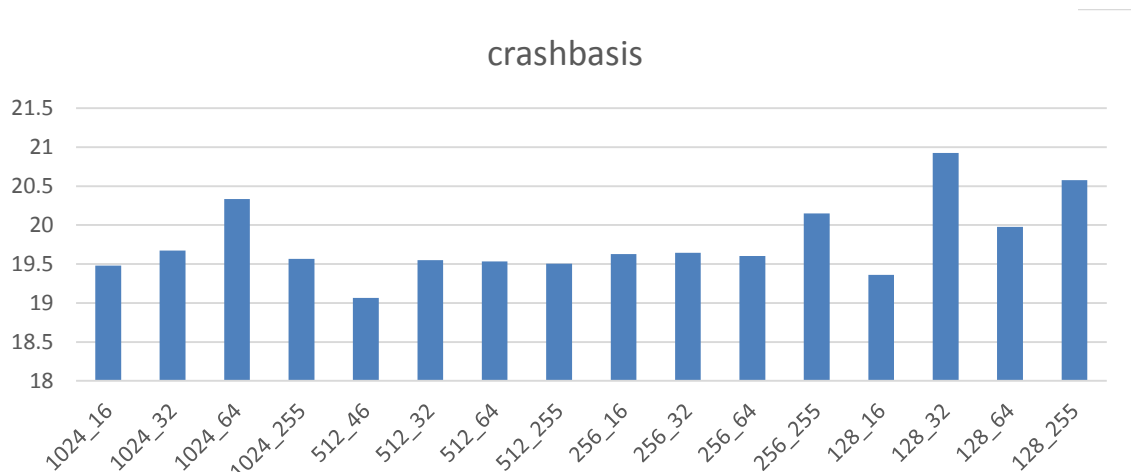
شکل ۴. مقدار کارایی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس تنک نمونه در قالب VCSR در اندازه‌های بلوک نخی ۱۲۸، ۲۵۶، ۵۱۲ و ۱۰۲۴



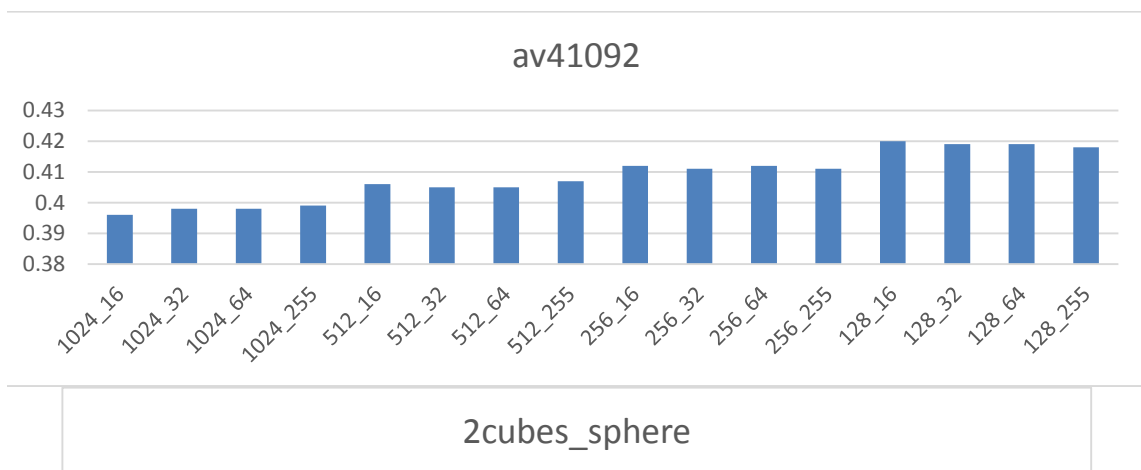
شکل ۶. مقدار کارایی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس تنک نمونه در قالب BCSR در اندازه‌های بلوک نخی ۱۲۸، ۲۵۶، ۵۱۲ و ۱۰۲۴



شکل ۷. مقدار کارآیی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس تنگ نمونه 2cubes\_sphere در قالب BELL با تعداد ثابت‌های ۱۶، ۳۲، ۶۴ و ۲۵۵

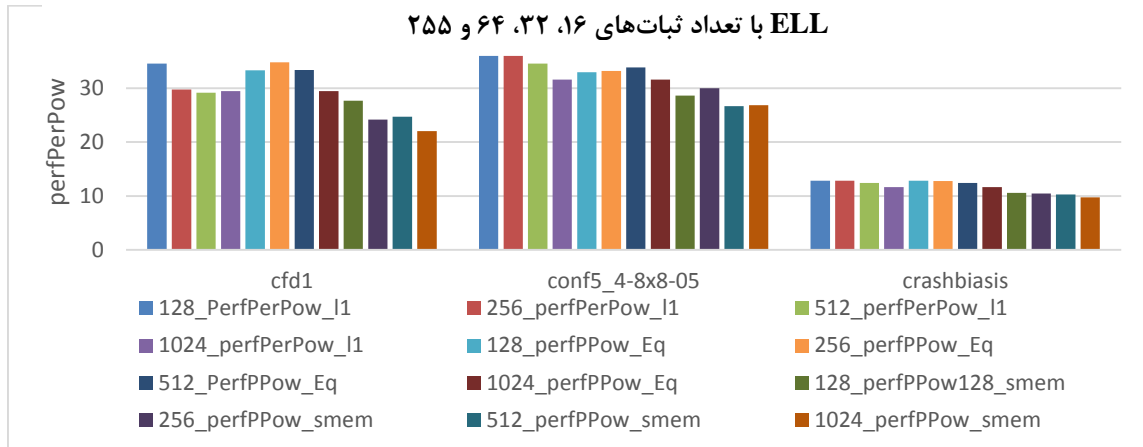


شکل ۸. مقدار کارآیی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس تنگ نمونه crashbasis در قالب CDS با تعداد ثابت‌های ۱۶، ۳۲، ۶۴ و ۲۵۵



2cubes\_sphere

شکل ۹. مقدار کارآیی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس تنک نمونه av41092 در قالب



شکل ۱۰. مقدار کارآیی در واحد توان مصرفی الگوریتم ضرب برداری برای ماتریس‌های تنک نمونه در قالب VCSR

**با تنظیمات مختلف سلسله مراتب حافظه**

مشخص است که با کاهش اندازه‌ی حافظه‌ی نهان، نمودار کارآیی در وات افت پیدا می‌کند.

**۶- نتیجه گیری**

همانگونه که پیشتر ذکر شد، مطالعات بر روی محاسبات تنک در قالب موارد بر روی کارآیی هسته‌ی ضرب برداری ماتریس‌ها تنک بر روی پردازنده‌ی گرافیکی تمرکز کرده‌اند و مطالعات انگشت شماری بر روی انرژی مصرفی این هسته صورت گرفته است. همچنین هیچ پژوهشی جهت طراحی سیستمی بهینه در مصرف انرژی با در نظر گرفتن کارآیی این هسته بر روی پردازنده‌ی گرافیکی انجام نشده است و امروزه که محققان بر روی لزوم در نظر گرفتن انرژی مصرفی در کنار کارآیی تاکید زیادی می‌کنند، برای عملیات ضرب برداری ماتریس‌های تنک که در اکثر برنامه‌های علمی ریاضی و فیزیک در مقیاس گسترده کاربرد دارند، نیاز به چنین مطالعه‌ای حس می‌شود. در این پژوهش مشاهده شد که رفتار الگوریتم ضرب برداری با توجه به قالبی که برای ذخیره سازی ماتریس تنک استفاده می‌کند متفاوت است و نمی‌توان یک پیکربندی را برای الگوریتم‌هایی با رفتارهای متفاوت در نظر گرفت. در این پژوهش تنظیمات مناسب برای سه پارامتر اندازه‌ی بلوک نخ، تعداد ثبات در نخ و سلسله مراتب حافظه برای مشهورترین و بهینه ترین قالب‌های ذخیره سازی تنک به دست آمده است.

نتایج آزمایشات با تنظیمات مختلف ۳ پارامتر ذکر شده برای ۲۰۰ ماتریس نمونه نشان می‌دهند که رفتار الگوریتم با استفاده از هر قالب در تمام آزمایشات یکسان است چرا که در همه‌ی موارد همیشه یکی از ترکیبات این تنظیمات بهترین کارآیی در وات را نتیجه می‌دهد. جدول ۱ تنظیمات مناسب برای هر قالب را که از آزمایشات بر روی ۲۰۰ ماتریس نمونه‌ی تنک به دست آمده، نشان می‌دهد.

**هر قالب**

نام قالب	اندازه‌ی بلوک نخ	تعداد ثبات در نخ	اندازه‌ی حافظه‌ی نهان (KB)
BCSR	256	16	48
CSR	256	16	48
SELL	256	32	48
ELL	128	16	48
BELL	128	32	48
VCSR	128	32	48
HYB	128	16	48
CDS	128	16	48
DIA	128	32	48

با این‌که قالب VCSR از حافظه‌ی اشتراکی برای ذخیره‌ی داده‌های متناظر با هر بلوک نخ استفاده می‌کند اما در شکل ۳-۲۸ مشاهده می‌شود که برای ماتریس‌های مختلف در این قالب کارآیی در وات با پیکربندی حافظه با بیشترین اندازه‌ی حافظه‌ی نهان به حداکثر خود رسیده است. در این نمودارها

Available on line at:  
<http://developer.nvidia.com/cuda-toolkit-80>.

11. Benatia, W Ji, Y Wang, F Shi, "Energy evaluation of Sparse Matrix-Vector Multiplication on GPU" in Green and Sustainable Computing Conference ,2016, p. 1-6. N.

12. Bell and M. Garland, "Implementing sparse matrix-vector multiplication throughput-oriented processors," In *Proc. of Int'l Conf. on High Performance Computing Networking, Storage and Analysis, SC '09*, pages 18:1-18:11. ACM, 2009.

13. S. Mullen, M. M. Wolf, and A. Klein, "Pakck: Performance and power analysis of key computational kernels on cpus and gpus," in *High Performance Extreme Computing Conference (HPEC), 2013 IEEE*, 2013, pp. 1-6.

14. H. Anzt, S. Tomov, and J. Dongarra, "Energy efficiency and performance frontiers for sparse computations on GPU supercomputers," in *Proceedings of the sixth international workshop on programming models and applications for multicores and manycores*, 2015, pp. 1-10.

15. Burtscher, I. Zecena, and Z. Zong, "Measuring GPU power with the K20 built-in sensor," in *Proceedings of Workshop on General Purpose Processing Using GPUs*, 2014, p. 28.

16. S. Song, C. Su, B. Rountree, and K. W. Cameron, "A simplified and accurate model of power-performance efficiency on emergent gpu architectures," in *Parallel & Distributed Processing (IPDPS), 2013 IEEE 27th International Symposium on*, 2013, pp. 673-686.

17. Zardoshti, P., Khunjush, F. & Sarbazi-Azad, H, "Adaptive sparse matrix representation for efficient matrix-vector multiplication," *Journal of Supercomputing* (2016) 72: 3366.  
<https://doi.org/10.1007/s11227-015-1571-0>.

18. NVIDIA, "NVIDIA Corporation (2014) Tuning CUDA applications for Kepler," *Technical report*, August 2014.  
[http://docs.nvidia.com/cuda/pdf/Kepler\\_Tuning\\_Guide.pdf](http://docs.nvidia.com/cuda/pdf/Kepler_Tuning_Guide.pdf)

## منابع

1. J. Im and K. Yelick, "Optimization of Sparse Matrix Kernels for Data Mining," in *Proc. of the Workshop on Text Mining*, 2001.
2. L. N. Trefethen and D. Bau, III, *Numerical Linear Algebra*. Society for Industrial and Applied Mathematics, 1997.
3. R. Gilbert, S. Reinhardt, and V. B. Shah, "Highperformance Graph Algorithms from Parallel Sparse Matrices," in *Proc. of the Int'l Workshop on Applied Parallel Computing*, 2006.
4. Owens JD, Luebke D, Govindaraju N, Harris M, Krüger J, Lefohn AE, Purcell TJ (2007) "A survey of general-purpose computation on graphics hardware," In: *Computer graphics forum*, vol 26. Wiley Online Library, pp 80–113.
5. Bell and M. Garland, "Efficient sparse matrix-vector multiplication on CUDA," *Nvidia Technical Report NVR-2008-004*, Nvidia Corporation 2008.
6. S. Yan, C. Li, Y. Zhang, and H. Zhou, "yaspmv: Yet another spmv framework on gpus," in *ACM SIGPLAN Notices*, 2014, pp. 107-118.
7. J. W. Choi, A. Singh, and R. W. Vuduc, "Model-driven autotuning of sparse matrix-vector multiply on GPUs," in *ACM Sigplan Notices*, 2010, pp. 115-126.
8. Bolz, Ian Farmer, Eitan Grinspun, and Peter Schroeder, "Sparse matrix solvers on the GPU: Conjugate gradients and multigrid," *ACM Trans. Graph.*, 22(3):917-924, 2003.
9. R. Gilbert, S. Reinhardt, and V. B. Shah, "Highperformance Graph Algorithms from Parallel Sparse Matrices," in *Proc. of the Int'l Workshop on Applied Parallel Computing*, 2006.
10. NVIDIA, "NVIDIA Corporation. *CUDA Toolkit Reference Manual*, 8.0 edition.