



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده مهندسی برق

گزارش سمینار در میکرو دو
گرایش الکترونیک دیجیتال

معماری پردازنده گرافیکی پاسکال

نگارش
علی زمانی

استاد درس
دکتر سید احمد معتمدی

بهمن ماه 1397

1	فصل اول مقدمه.....
3	فصل دوم معماری پردازنده‌های گرافیکی.....
4	2.1 کارتهای گرافیک.....
7	2.2 معماری CUDA.....
12	2.3 حافظه در کودا.....
15	2.4 معماری پردازنده های گرافیکی Geforce سری 8.....
24	2.5 معماری پردازنده های گرافیکی GeForce سری 200.....
27	2.6 معماری Fermi.....
31	2.6.1 SMهای نسل سوم.....
32	2.6.2 حافظه با قابلیت بازپیکربندی.....
34	2.6.3 واحد تصحیح و تشخیص خطا.....
35	2.6.4 تکنولوژی PTX 2.0.....
36	2.6.5 اجرای همزمان کرنل ها.....
37	2.6.6 کاستیها فرمی.....
39	2.6.7 NVIDIA TESLA C2050.....
41	3. فصل سوم معماری پاسکال.....
42	3.1 Tesla P100: کارایی بی نظیر در کنار ویژگی‌های بارز و چشم نواز برای پردازش های گرافیکی.....
43	3.2 Tesla P100 پیشنهادی عالی برای تکنولوژی NVLink و سرورهای PCI-Express.....
44	3.3 معماری پاسکال GP100: از هر جهتی سریعتر.....
46	3.4 Pascal SM.....
48	3.5 عملکرد بالا در دقت مضاعف.....
49	3.6 افزایش سرعت عملکرد عملگر های Atomic.....
50	3.7 بهبودی شگرف در حافظه.....
53	3.8 ECC Memory.....
53	3.9 NVLink برای ارتباط سریع.....
56	3.10 افزایش بهره وری برنامه نویسی با استفاده از حافظه یکپارچه.....
57	3.10.1 حافظه یکپارچه پردازنده P100.....
59	3.11 NVidia DGX-1 یادگیری ژرف با استفاده از سوپر کامپیوترها.....
	4. فصل چهارم تسریع الگوریتم های یادگیری ژرف و هوش مصنوعی با استفاده از پردازنده
61	گرافیکی.....

4.1	یادگیری ژرف در نگاه کلی.....	63
4.2	پردازنده های گرافیکی NVidia: موتوری قوی برای یادگیری ژرف.....	67
4.3	نرم افزار ها و toolkit های لازم برای یادگیری ژرف.....	69
71	فصل چهارم جمع بندی و نتیجه گیری و پیشنهادات جمع بندی و نتیجه گیری.....	
73	منابع و مراجع.....	

شکل 1: ترد، ترد بلاک و حافظه‌ی سراسری در پردازنده‌های گرافیکی.....	9
شکل 2: کرنل، گرید، بلاک و ترد در مدل میزبان - وسیله.....	10
شکل 3: مدل میزبان، وسیله.....	11
شکل 4: بانکهای حافظه در کودا.....	14
شکل 5: معماری پردازنده‌ی گرافیکی GeForce 8800.....	16
شکل 6: TPC.....	19
شکل 7: نمونه‌ی تار پارچه.....	21
شکل 8: اجرای SIMT.....	21
شکل 9: معماری GeForce GT 200.....	25
شکل 10: اجزای واحد TPC در پردازنده‌های گرافیکی سری 200.....	26
شکل 11: معماری فرمی.....	31
شکل 12: ساختار حافظه در فرمی.....	34
شکل 13: فضای آدرس یکتا در معماری فرمی.....	36
شکل 14: اجرای هم‌زمان کرنلها در فرمی.....	37
شکل 15: پردازنده‌ی تسلا P100.....	43
شکل 16: بلوک دیگرامی از P100.....	44
شکل 17: SMهای پردازنده‌ی تسلا P100.....	47
شکل 18: اشتابدهنده p100 از جلو.....	Error! Bookmark not defined.
شکل 19: اشتابدهنده‌ی P100 از پشت.....	51
شکل 20: نمایشی از شمای دای و میکروباپ ها.....	52
شکل 21: محل اینترپوزر.....	52
شکل 22: NVLink.....	55
شکل 23: هشت پردازنده گرافیکی که به صورت hyper cube قرار دارند.....	55
شکل 24: 4 پردازنده گرافیکی که با استفاده از تکنولوژی NVLink به پردازنده اصلی و خود متصل اند.....	56
شکل 25: نمونه کد کودا 6.....	57
شکل 26: یک سوپر کامپیوتر DGX-1.....	60
شکل 27: یک پرسپترون که مدلی ساده از شبکه عصبی است.....	64
شکل 28: یک مدل پیچیده از شبکه عصبی که نیازمند محاسبات سنگین است.....	66
شکل 29: فریم وورک های مختلف برای کار با یادگیری ژرف.....	70

فصل اول

مقدمه

مقدمه

پردازنده های گرافیکی امروزه در زمینه های گوناگونی از قبیل بازی ها، نرم افزار های گرافیکی همچون فوتوشاپ ، corel و ... مورد مصرف قرار میگیرند. با روی کار آمدن الگوریتم های شبکه ژرف و هوش مصنوعی در چندین سال اخیر و نیاز محاسباتی زیاد آنها، پردازنده های گرافیکی راه خود را در این بین باز کرده و سر دم دار عرصه سخت افزار های قدرتمند برای پردازش های موازی در حیطه شبکه های چند لایه ژرف شده اند.

شرکت NVidia با عرضه معماری جدیدی تحت عنوان پاسکال این امکان را به برنامه نویسان و محققان داده است که برنامه های شبکه عصبی و یادگیری ژرف خود را در کمتر از دقیقه به اجرا در آورند، الگوریتم هایی که اجرای آنها بر روی سخت افزار های سنتی ممکن بود تا روز ها به طول انجامد. این شرکت با همکاری های وسیع با کمپانی های نرم افزاری و اینترنتی دنیا از قبیل فیسبوک و نتفلیکس و گوگل مسیر را برای بسیاری از الگوریتم های یادگیری ژرف هموار ساخته است. همچنین بسیاری از کمپانی های ماشین سازی دنیا نظیر آوودی و بنز با بهره گیری از این پردازنده های گرافیکی انقلابی در پدیده ماشین های بدون راننده کرده اند.

فصل دوم

معماری پردازنده های گرافیکی

پردازنده گرافیکی در یک نگاه

انسان ها قادر به درک سه بعد هستند و از دنیای اطرافشان دیدی سه بعدی دارند. همین ویژگی انسان ها سبب شده است که از سه دهه پیش تاکنون تلاش های بسیاری برای شبیه سازی دنیای سه بعدی خارجی بر روی صفحه ی نمایشگر رایانه ها که دارای دو بعد است صورت پذیرد. این تلاش ها که متمرکز بر نمایش سه بعدی یک تصویر دوبعدی است، به محاسبات بسیار سنگینی نیاز دارد این محاسبات در صورتی که قرار با شد با کیفیت بالایی انجام شوند، از عهده ی CPU¹ ها که پردازنده هایی قدرتمند و همه منظوره هستند نیز خارج است. البته ممکن است با توجه به توان CPU و نوع کاری که مدنظر است، امکان واگذار کردن پردازش گرافیکی به CPU ممکن باشد اما این کار به قیمت بازداشتن CPU از رسیدن به کارهای اصلی اش تمام خواهد شد. به همین دلایل با مرور زمان تصمیم بر این شده است که پردازنده ای جداگانه برای پردازش های گرافیکی اختصاص پیدا کند. برخی از افراد این پردازنده، یعنی واحد پردازش گرافیکی² را با کارت گرافیک³ یکی می دانند درحالی که بین این دو قسمت تفاوت های بارزی وجود دارد. کارت گرافیک قطعه ای است که به کمک آن رایانه قادر می گردد تصویر را تولید کرده و برای نمایش در اختیار نمایشگر قرار دهد. یکی از قسمت های کارت گرافیک واحد پردازش گرافیکی است. در این گزارش به واحد پردازش گرافیکی به اختصار پردازنده ی گرافیکی گفته می شود. پردازنده ی گرافیکی مهم ترین قسمت یک کارت گرافیک است که مسئول اجرای پردازش های لازم برای نمایش تصویر است. در این فصل ما به صورت خلاصه در مورد کارت های گرافیک توضیح خواهیم داد، سپس بحث پردازش موازی با کمک پردازنده های گرافیکی را با توضیح معماری CUDA⁴ آغاز می کنیم و سپس بحث را با شرح معماری پردازنده های گرافیکی ساخت شرکت NVIDIA ادامه خواهیم داد.

2.1 کارت های گرافیک

نخستین کارت گرافیک در سال 1960 میلادی توسط شرکت IBM ساخته شد. این کارت گرافیک که همراه رایانه های IBM PC فروخته می شد دارای 4 کیلوبایت حافظه بود و فقط قادر بود رنگ سبز را نشان دهد که از آن برای نمایش کارکترها استفاده می شد. هرچند نخستین کارت گرافیک دارای قدرت پردازشی بالایی نبود اما همین کارت ساده نویدبخش تحول بزرگی در صنعت پردازش الکترونیک بود. با توجه به اینکه تصویر دارای جذابیت های خاص خودش است، هر تغییری که در ساختار کارت های گرافیک ایجاد و موجب بهبود تصویر نهایی می شد از سوی کاربران با استقبال فراوان روبه رو می شد. از طرف دیگر با رشد تدریجی کارت های گرافیکی و

¹Central Processing Unit

² Graphics Processing Unit

³Graphic Card

⁴Compute Unified Device Architecture

افزایش قدرت آن ها صنعت بازی های رایانه ای شروع به کار کرد. همچنین وجود کاربران متعدد باعث تضمین سود ناشی از فعالیت در بازار تو سعه کارت های گرافیکی شد. بدین سان صنعت تو سعه کارت های گرافیکی با سرعتی فزاینده شروع به رشد کرد. وظیفه ی کارت گرافیک این است که سیگنال های دیجیتال تولید شده توسط پردازنده را دریافت کرده و آن ها را تبدیل به سیگنال های آنالوگ یا دیجیتال قابل نمایش توسط نمایشگر بکند. شایان ذکر است که این کار باید در زمان بسیار کمی انجام شود به گونه ای که تغییرات برای چشم انسان آزاردهنده نباشد.

هر کارت گرافیک دارای واحدهای مختلفی است که هر کدام برای کار خاصی در نظر گرفته شده اند. این واحدها در کارت های مختلف ممکن است تغییر کنند. از آنجایی که بحث در مورد طراحی و جزئیات کارت های گرافیکی از موضوع بحث این گزارش خارج است، در ادامه برخی از قسمت های مهم به صورت خلاصه توضیح داده خواهند شد.

پردازنده ی گرافیکی: این واحد مهم ترین قسمت کارت گرافیک است که عملاً مسئولیت پردازش تصاویر را بر عهده دارد. بحث ما در این کتاب روی این واحد متمرکز است و توضیح خواهیم داد که چگونه می توان از این واحد برای پردازش موازی استفاده کرد.

حافظه: حافظه برای ذخیره سازی اطلاعات مربوط به تصاویر و کار بر روی آن هاست. در کارت های گرافیک جدیدتر از حافظه هایی با ظرفیت زیاد استفاده می شود زیرا در بسیاری از مسائل پردازش موازی با ابعاد وسیع حجم حافظه تأثیر به سزایی بر روی سرعت حل مسئله خواهد داشت.

ادوات جانبی: در هر کارت گرافیک یک سری ادوات جانبی نیز وجود دارد. برای مثال می توان به واسط ارتباط با برد اصلی¹ اشاره کرد؛ این واسط می تواند PCI²، PCIe³، AGP⁴ و ... باشد در حال حاضر⁵ بیشتر پردازنده های گرافیکی از PCI e استفاده می کنند. درگاه های⁶ خروجی تصاویر نیز از دیگر ادوات جانبی هستند. این درگاه ها که برای انتقال اطلاعات بر روی نمایشگر ساخته شده اند ممکن است دارای استاندارد DIV⁷ یا VGA⁸ باشند. کارت های گرافیکی دارای یک قسمت بایاس هم هستند که عملاً برنامه ای است که کارت گرافیک به کمک آن با CPU ارتباط برقرار می کند و دستورهای CPU را اجرا می کند. از آنجایی که پردازنده های گرافیکی گرمای زیادی تولید می کند کارت های گرافیکی معمولاً دارای چندین مبدل حرارتی و پنکه برای تعدیل دما هستند. در هنگام نصب

¹Mainboard

²Peripheral Component Interconnect

³PCI Express

⁴Accelerated Graphic Port

سال 2013

⁵Port

⁶Digital Video Interface

⁷Video Graphics Array

کارت گرافیک نیز باید مسائل مربوط به خنک شدن در نظر گرفته شود. از دیگر مزیت های موجود در کارت های گرافیک می توان به خطوط لوله بازنمایی^۱، واحدهای بافت و سایه زنی اشاره کرد.

همان طور که پیش تر ذکر شد با وارد شدن بازی های رایانه ای سرعت توسعه کارت های گرافیک بیشتر شد زیرا با ارائه هر کارت گرافیک قدرتمندتر بازی های جذاب تری قابل ساخت بود و هر بازی جدید خود پردازنده ی گرافیکی قوی تری را طلب می کرد. این حلقه ی تشدید شونده و علاقه ی روزافزون قشر وسیعی از کاربران به بازی های رایانه ای از یک سو و جذابیت های بصری رایانه برای تمامی کاربران از سوی دیگر باعث شد که روز به روز پردازنده های گرافیکی قدرتمندتری تولید شود و کاربران به راحتی آن ها را خریداری کنند. در این بین برخی از کاربران به این فکر افتادند که از قدرت بالای پردازنده های گرافیکی برای کارهای مفیدتری! استفاده کنند. این نقطه ی آغازی بر استفاده از پردازنده های گرافیکی برای کارهای عمومی و علمی بود. واقعیت این است که در حال حاضر بسیاری از رایانه ها دارای پردازنده های گرافیکی بسیار قدرتمندی هستند. این پردازنده ها تا زمانی که رایانه درگیر انجام کارهای پردازش تصاویر پیچیده است توانایی خارق العاده ای به آن می دهند اما به صورت کلی اکثر رایانه ها وقت بسیار کمی را برای پردازش تصاویر اختصاص می دهند و در سایر زمان ها پردازنده گرافیکی بیکار مانده، انرژی مصرف کرده و آن را تبدیل به گرما می کند. با در نظر گرفتن این نکات استفاده از پردازنده گرافیکی برای سایر مقاصد پردازشی از هر جنبه مفید به نظر می رسد هر چند که این کار در ابتدا کار دشواری بود زیرا پردازنده های گرافیکی اولیه دارای دو واحد پردازش پیکسلی^۲ و پردازش رأسی^۳ بودند. واحد پردازش رأسی مسئول انجام محاسبات رأسی بود. برای مثال در واحد پردازش رأسی باید مختصات یک تصویر سه بعدی در یک فضای دوبعدی (صفحه ی مانیتور) به دست آید. واحد پردازش پیکسلی هم مسئول مشخص کردن رنگ بخش های مختلف تصویر بود. با توجه به این نکات در آغاز کار برای پردازش موازی با پردازنده ی گرافیکی می بایست برنامه ها به فرم مسائل پردازش تصویر در بیابند و سپس برای پردازش به پردازنده ی گرافیکی داده شوند. این کار نه تنها دشوار بود بلکه برای تمامی مسائل هم عملی نبود. به تدریج شرکت های سازنده ی پردازنده های گرافیکی سعی کردند که پردازنده ی خاص منظوره ی خود را از خاص منظور بودن درآورده و آن را برای اجرای عملیات پردازش موازی بهینه نمایند. برای نیل به این هدف برداشتن دو گام لازم بود. گام نخست تغییر معماری سخت افزار و ارائه معماری مناسب پردازش

^۱ Rendering pipeline

^۲ Pixel Fragment Processor

^۳ Vertex Processor

موازی بود و گام دیگر توسعه ی زبان های برنامه نویسی مناسب برای فراهم کردن امکان کار با پردازنده ی گرافیکی بود. شرکت NVIDIA این دو گام را برداشت و یک معماری تحت عنوان ¹CUDA به وجود آورد.²

2.2 معماری ³CUDA

کودا که نخستین بار در سال 2006 هم زمان با ارائه ی GeForce 8800 رونمایی شد، یک بستر نرم افزاری⁴ برای توسعه ی برنامه های موازی بر روی برخی⁵ پردازنده های گرافیکی ساخت شرکت NVIDIA را فراهم می کند. با استفاده از این بستر نرم افزاری، می توان روی پردازنده های گرافیکی دقیقاً مانند CPU برنامه نویسی کرد. به این استفاده از پردازنده های گرافیکی برای حل مسائل عمومی ⁶GPGPU گفته شده که گاهی اوقات به اختصار به صورت GPU² نوشته می شود. امروزه NVIDIA به جای GPGPU از اصطلاح GPU Computing استفاده می کند. کودا به صورت یک افزونه برای زبان های برنامه نویسی استاندارد C، C++ و Fortran وجود دارد. شایان ذکر است که این زبان ها تنها راه استفاده از پردازنده های گرافیکی نیست بلکه امکان کار با پردازنده های گرافیکی از طرق دیگری نیز مانند استفاده از OpenCL و یا استفاده از جعبه ابزارهای متلب و LabVIEW میسر است. مدلی که کودا برای پردازش ارائه می کند به صورت زیر است.

1. حافظه ی مورد نیاز بر روی حافظه ی پردازنده ی گرافیکی رزرو می شود.
2. اطلاعات از حافظه ی اصلی یا ورودی ها برداشته شده و به حافظه ی پردازنده ی گرافیکی انتقال داده می شود.
3. CPU دستوراتی را که باید اجرا شوند برای اجرا به پردازنده ی گرافیکی می فرستد.⁷
4. پردازنده ی گرافیکی دستورات را به صورت موازی بر روی داده هایی که در مرحله ی یک دریافت کرده است اجرا می کند.

¹Compute Unified Device Architecture

² شایان ذکر است که هر چند CUDA مخفف عبارت Compute Unified Device Architecture است اما این عبارت مخفف فقط قبلاً به کار گرفته می شد. در حال حاضر شرکت NVIDIA این عبارت مخفف را حذف کرده است و فقط از CUDA به عنوان یک نشانه (Brand) استفاده می کند، نشانه ای که امروزه (2013) خیلی معروف شده است.

³ برای راحتی از اینجا به بعد واژه ی CUDA را به صورت «کودا» خواهیم نوشت.

⁴Platform

⁵ پردازنده های گرافیکی قدیمی از این بستر پشتیبانی نمی کنند. برای اینکه بدانید که آیا یک پردازنده ی گرافیکی از این بستر پشتیبانی می کند یا

خیر به دفترچه راهنما یا سایت شرکت NVIDIA مراجعه نمایید.

⁶General-Purpose computing on Graphics Processing Units

⁷ Run kernel

5. اطلاعات پس از همزمان سازی CPU با پردازنده ی گرافیکی از حافظه ی پردازنده ی گرافیکی برداشته شده و به مقصد (حافظه ی اصلی یا خروجی) انتقال پیدا می کند.

6. حافظه ای که در پردازنده ی گرافیکی رزرو شده است آزاد می گردد.

هرچند که پیش از کودا نیز امکان برنامه نویسی موازی برای مقاصد عمومی بر روی پردازنده های گرافیکی وجود داشت، اما مدلی که کودا ارائه می کند، نسبت به مدل های پیش تر دارای چند ویژگی چشم گیر است:

- توانایی دستیابی به آدرس های مختلف حافظه برحسب نیاز برنامه نویس بدون وجود محدودیت جدی سخت افزاری؛ توانایی دستیابی اتفاقی به حافظه.
- در نظر گرفته شدن حافظه ی مشترک در بلاک ها. این حافظه ی مشترک امکان تبادل داده بین بلاک ها را با پهنای باند بیشتر و تأخیر کمتری فراهم می کند.
- پشتیبانی کامل از عملیات ریاضی بر روی اعداد صحیح¹ و پشتیبانی از عملیات منطقی بیتی.
- همزمان سازی منطقی و آسان تردها
- امکان فراگیری آسان

یک برنامه ی CUDA متشکل از یک سری کرنل² است که ممکن است موازی باهم یا به صورت متوالی اجرا شوند. هر کرنل به صورت موازی در قالب چندین ترد اجرا می گردد. برای راحت تر شدن برنامه نویسی این تردها به صورت یک سری دسته که به آن ها بلاک³ گفته می شود تقسیم بندی می شوند. خود بلاک ها هم در قالب یک سری گرید⁴ مرتب می شوند. پس هر گرید شامل تعدادی بلاک و هر بلاک شامل تعدادی ترد است. پردازنده ی گرافیکی برنامه ی کرنل را بر روی تعدادی ترد به صورت همزمان اجرا می کند. هر ترد به تنهایی کل برنامه ی کرنل را در اختیار دارد اما با تدابیری که بعداً توضیح داده خواهد شد، هر ترد فقط بخشی از برنامه ی کرنل را اجرا می کند بدین ترتیب زمانی که کار کل تردها به پایان برسد عملاً اجرای کرنل به پایان رسیده است. هر ترد دارای یک شماره ی اختصاصی، یک شمارنده ی برنامه⁵، تعدادی رجیستر و حافظه ی اختصاصی انحصاری است. در شکل 0-1 این تقسیم بندی ها نشان داده شده است.

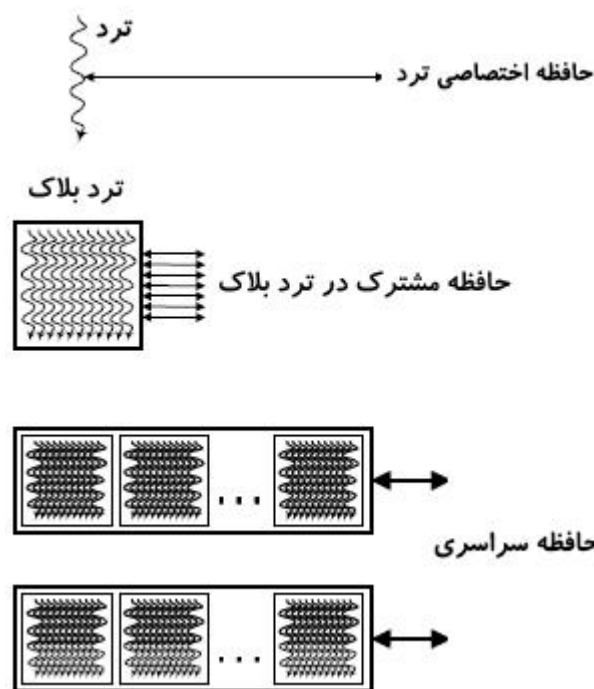
¹ نمونه های جدید پردازنده های گرافیکی منطبق با کودا از اعداد اعشاری نیز به خوبی پشتیبانی می کنند.

²Kernel

³Block

⁴Grid

⁵Program Counter



شکل 0-1: ترد، ترد بلاک و حافظه ی سراسری در پردازنده های گرافیکی

در معماری کودا ساختمان چینش ترد بلاک ها در داخل گرید می تواند یک بعدی یا دوبعدی باشد. درعین حال ساختمان تردها هم در داخل ترد بلاک می تواند یک، دو یا سه بعدی باشد. در فصول آتی در مورد ترد، ترد بلاک، گرید و کرنل بحث خواهیم کرد. در حال حاضر فقط کافی است با این مفاهیم یک آشنایی سطحی داشته باشید برای اینکه بحث به روشنی ختم گردد اجازه دهید مفاهیم فوق الذکر را مجدداً مرور کنیم.

کرنل: به قسمت یا قسمت های کوچکی از برنامه که برای اجرا بر روی پردازنده ی گرافیکی ارسال می گردد کرنل می گویند. کرنل قسمتی از برنامه بوده که دارای ویژگی های زیر است.

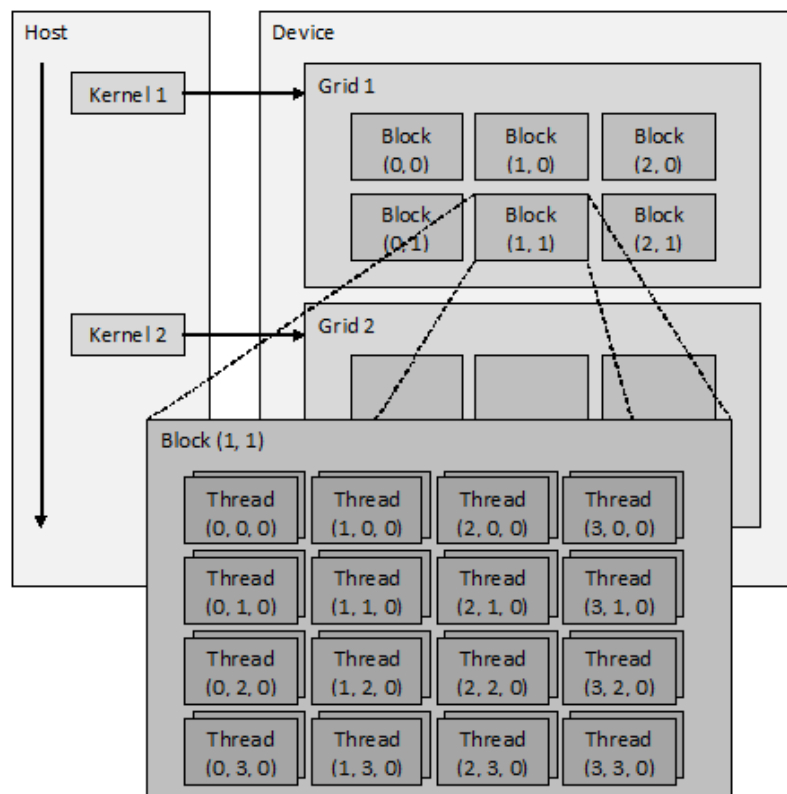
- بتوان به آن به صورت یک تابع ایزوله شده نگاه کرد.
- قادر باشد بر روی داده های مستقل به صورت مستقل عمل کند.

ترد: به کوچک ترین قسمت برنامه که برای اجرای موازی در نظر گرفته شده است ترد می گویند. یک Thread از تعدادی دستورالعمل تشکیل می شود که به ترتیب و بر روی یک هسته کودا اجرا می شود. علت انتخاب نام Thread (به معنای نخ) نیز به همین علت است. درواقع، Thread کوچک ترین جزء ساختار سلسله مراتبی مدل

کودا است که اجراکننده ی دستورالعمل های کرنل است. یعنی درنهایت، یک کرنل توسط تعداد زیادی ترد اجرا می شود.

ترد بلاک: در داخل هر گرید تعدادی ترد بلاک (بلاک) وجود دارد. تردهایی که در داخل ترد بلاک هستند در حین اجرای برنامه می توانند باهم داده رد و بدل نمایند.

گرید: در داخل هر کرنل فقط یک گرید وجود دارد. گرید یک تعریف است که نحوه ی اختصاص منابع به کرنل را نشان می دهد. به مجموعه ای از ترد بلاک ها یک گرید می گویند. در شکل 0-2 رابطه ی کرنل و گرید نشان داده شده است.

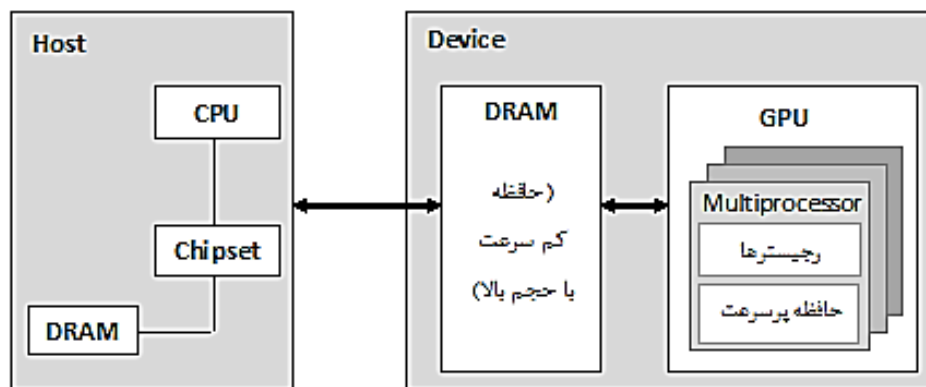


شکل 0-2: کرنل، گرید، بلاک و ترد در مدل میزبان - وسیله

پیش تر گفتیم که پردازنده های قدیمی تر دارای دو واحد مهم بودند که عبارت اند از واحد پردازش رأسی و واحد پردازش پیکسلی که به ترتیب مسئول اجرای عملیات سایه زنی روی رئوس و پیکسل ها هستند. عملیات سایه زنی شامل تمامی کارهایی است که برای ایجاد سطوح مناسبی از نور بر روی یک تصویر صورت می گیرد. عملیات سایه

زنی شامل دو قسمت مهم است. قسمت نخست سایه زنی رئوس^۱ است. در این قسمت یکی از وظایف مهم یافتن محل رئوس سه بعدی در یک مختصات دوبعدی است. قسمت دیگر که با نام واحد سایه زنی پیکسل^۲ شناخته می شود نیز برای محاسبه ی مقادیر رنگ و سایر مشخصات هر بخش از تصویر به کار گرفته می شود. متأسفانه این معماری برای پردازش موازی خیلی مناسب نبود به همین دلیل شرکت NVIDIA سعی کرد معماری پردازنده های گرافیکی را به نحوی بهینه کند تا هم امکان اجرای برنامه های گرافیکی را به همان صورت قبل داشته باشند هم اینکه بتوان از آن ها برای پردازش موازی استفاده شود. اولین تلاش شرکت NVIDIA در این راه منجر به تولید پردازنده ی گرافیکی GeForce سری 8 شد.

برای اینکه صحبت در مورد برنامه نویسی موازی ناهمگون با مشارکت همزمان CPU و GPU آسان تر شود، در بسیاری از مواقع از تشریح روند کار به کمک مدل میزبان^۳، وسیله^۴ استفاده می شود. این مدل در شکل 0-3 نشان داده شده است. همان طور که در شکل نشان داده شده است، پردازنده ی اصلی قادر است فقط از طریق حافظه ی خارجی به تبادل داده با پردازنده ی گرافیکی بپردازد. لازم به ذکر است که هرچند پیاده سازی واقعی با این مدل خیلی تفاوت دارد اما این مدل برای به دست آوردن یک درک کلی از روش کار مناسب است. دقت کنید که اولاً سرعت دسترسی به منابع داخلی موجود پردازنده ی گرافیکی بسیار بیشتر از سرعت دسترسی به منابع مشترک است. ثانیاً پردازنده ی اصلی نمی تواند به منابع داخلی پردازنده ی گرافیکی دسترسی پیدا کند و بالعکس.



شکل 0-3: مدل میزبان، وسیله

¹Vertex Shading

²Pixel Shading

³ Host

⁴ Device

2.3 حافظه در کودا

در یک پردازنده ی گرافیکی منطبق بر کودا، شش نوع حافظه ی مختلف با کاربردهای متفاوت وجود دارد. این حافظه ها را می توان به دودسته ی حافظه های روی تراشه و حافظه های بیرون تراشه تقسیم کرد. حافظه های بیرون تراشه از نظر فیزیکی جزئی از یک DRAM محسوب شده و لذا ویژگی های DRAM را دارند یعنی دارای حجم بالا و سرعت پایین (در صورتی که cache نشوند، بین 200 الی 300 سیکل) می باشند. این حافظه ها شامل موارد زیر هستند.

- **حافظه سراسری^۱:** هر پردازنده ی گرافیکی یک نمونه از این حافظه در اختیار دارد. این حافظه از سمت پردازنده ی گرافیکی هم قابل خواندن و هم قابل نوشتن است. شایان ذکر است که این حافظه برخلاف حافظه های ثوابت و بافت دارای حافظه ی نهان نیست.
- **حافظه ثوابت^۲:** هر پردازنده ی گرافیکی یک نمونه از این حافظه را در اختیار دارد. این حافظه از سمت پردازنده ی گرافیکی فقط قابل خواندن است. حافظه ی ثوابت همان طور که از نامش پیداست، حافظه ای است برای ذخیره کردن مقادیر ثابت. این مقادیر، می توانند آرایه ای از اعداد باشند که قبلاً توسط پردازنده ی اصلی به این حافظه منتقل شده و ترد آن ها را می خوانند. طول عمر اطلاعات موجود در حافظه ثوابت برابر با طول عمر برنامه است و در طول مدتی که کرنل های مختلفی از یک برنامه اجرا می شوند، حافظه ثوابت مقادیر خود را حفظ می کند. این حافظه، دارای حافظه ی نهان نیز هست. حافظه ی نهان برای کاربردهایی که در آن ها پردازنده ی گرافیکی به صورت مکرر اعدادی را از حافظه ی خود می خواند، استفاده می گردد. سرعت دسترسی به اطلاعات موجود در حافظه نهان برابر با سرعت دسترسی به حافظه مشترک است پس به صورت کلی سرعت این حافظه می تواند بین 1 تا 300 سیکل ساعت تغییر کند. باید توجه داشت که استفاده ی اشتباه و یا بیش از حد نیاز از این حافظه تنها احتمال عدم اصابت^۳ را زیاد می کند.
- **حافظه بافت^۴:** هر پردازنده ی گرافیکی دارای یک حافظه ی بافت با دسترسی خواندن و نوشتن است. کلمه ی بافت در اصطلاح گرافیکی به معنای تصویری است که بر روی سطح اشیای سه بعدی نمایش داده می شود. در کودا این حافظه نیز کاربردهای منحصر به خود را دارد. این حافظه نیز دارای حافظه ی نهان است. با توجه به اینکه این حافظه، جزئی از DRAM به حساب می آید، سرعت دسترسی به اطلاعاتی از آن که در حافظه ی نهان موجود نیست بسیار پایین است. این حافظه برای حالتی بهینه سازی شده است که

¹ Global memory

² Constant memory

³ Cache miss

⁴ Texture memory

بخواهیم آرایه های دوبعدی ذخیره کنیم و تردهای بخواهند درایه هایی نزدیک به هم را از آن بخوانند. برای حافظه بافت، منطق ذخیره سازی اطلاعات در حافظه ی نهان مطابق با نیازهای پردازش گرافیکی بافتها است. علاوه بر این، حجم حافظه بافت قابل تنظیم است و حداکثر می تواند برابر با کل فضای DRAM باشد. اما استفاده ی بیش از حد از آن، تنها باعث کاهش نرخ اصابت^۱ و سرعت خواندن از آن می شود.

- **حافظه محلی^۲:** هر ترد بخشی از این حافظه را، به صورت انحصاری، برای خواندن و نوشتن در اختیار دارد. هرگاه اطلاعات یک ترد در رجیسترها جا نشود و نخواهیم این اطلاعات را در حافظه ی مشترک قرار دهیم باید این اطلاعات را به حافظه ی محلی سپرد. حافظه های محلی بر روی DRAM قرار دارند پس دارای حجم بسیار بزرگ تری نسبت به حافظه ی مشترک می باشند. هر ترد به بخشی از حافظه ی DRAM تحت عنوان حافظه ی محلی دسترسی دارد. شایان ذکر است که پس از پایان اجرای کرنل و از بین رفتن ترد، قسمتی از حافظه که در اختیار این ترد بود نیز از بین خواهد رفت یعنی اطلاعات آن دیگر قابل دسترسی نخواهد بود.

دسته ی دوم حافظه های روی تراشه هستند. این حافظه ها دارای سرعتی بالاتر (یک سیکل) و حجمی پایین تر (حافظه ی مشترک^۳ معمولاً 16 یا 48 کیلوبایت بوده و ثباتها معمولاً 16384 یا 8192 عدد هستند)^۴ هستند.

- **حافظه مشترک:** همه ی تردهای موجود در یک بلاک، به صورت اشتراکی از یک نمونه از این حافظه برای خواندن و نوشتن استفاده می کند.

- **ثباتها:** هر ترد می تواند از تعدادی ثبات داخلی برای خواندن و نوشتن استفاده کند.

چون تردهای موجود در یک بلاک باهم از حافظه ی مشترک استفاده می کنند، این حافظه نیاز به پهنای باند زیادی دارد که بتواند به همه ی این تقاضاها پاسخ دهد. برای دسترسی به این پهنای باند حافظه های مشترک به تعدادی بانک تقسیم شده اند. این بانکها دارای ورودی و خروجی مستقل و طول برابری هستند. برای مثال در پردازنده گرافیکی G80، هر حافظه مشترک 16 کیلوبایت بوده و از 16 بانک یک کیلوبایتی با کلمات 32 بیتی تشکیل می شود. هر بانک می تواند در هر سیکل به یک آدرس سرویس دهی کند. به عبارتی دیگر پهنای باند هر بانک، 32 بیت در هر سیکل ساعت است. بنابراین تعداد سرویس های همزمانی که حافظه مشترک می تواند بدهد، به تعداد بانکهایش است.

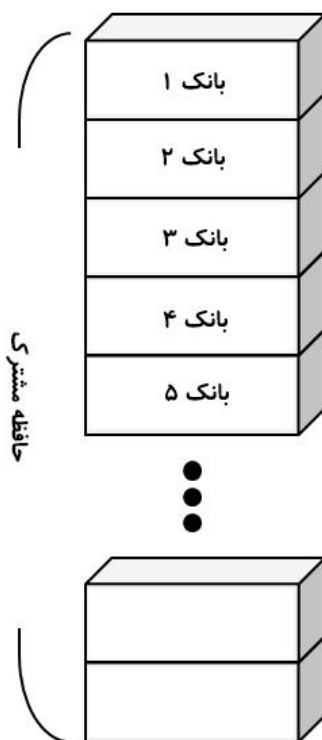
¹ Cache hit

² Local memory

³ Shared Memory

⁴ واضح است که این اعداد در آینده تغییر خواهند کرد.

شکل 4-0 نحوه تقسیم یک حافظه مشترک را به بانک های ذکر شده، نمایش می دهد. کلمات 32 بیتی متوالی، به بانک های متوالی تخصیص داده می شوند.



شکل 4-0: بانک های حافظه در کودا

اگر در آن واحد، چندین ترد تقاضای دسترسی به یک بانک را داشته باشند، تداخل پیش می آید و پردازنده ی گرافیکی ناچار است که به تقاضاها، به صورت نوبتی پاسخ دهد. در صورت بروز تداخل بانک، سرعت نهایی حافظه مشترک افت می کند. تقسیم حافظه مشترک به بانک های متعدد، احتمال تداخل در دسترسی به بانک ها را کاهش داده و سرعت نهایی خواندن و نوشتن در آن را افزایش می دهد.

در حالت عادی، تداخلی بین Thread های دونیمه ی Warp وجود ندارد؛ زیرا Thread های نیمه ی دوم Warp صبر می کنند تا Thread های نیمه ی اول کار شان با Shared Memory تمام شود. هرچند این اتفاق می تواند نوعی تداخل بانک محسوب شود، اما نظر به اینکه هر Warp شامل 32 Thread بوده و تنها 16 بانک مستقل برای پاسخ دهی به نیازهای این Thread ها وجود دارد، بنابراین تداخل بین دونیمه ی یک Warp امری است اجتناب ناپذیر و معمولاً از تأخیر به وجود آمده صرف نظر می شود.

با توجه به این که آدرس های یک Shared Memory معمولاً 32 بیتی هستند، در حالتی که با آرایه ای سروکار داریم که اعضای آن ها متغیرهای غیر 32 بیتی (مثل char یا int64) هستند، بایستی مراقب تداخل در بانک ها باشیم.

دو شکل 4-7، حالتی را معرفی می کنند که هیچ تداخلی در دسترسی Thread ها به بانک های Shared Memory وجود ندارد. شکل سمت چپ آدرس دهی خطی با گام یک، و شکل سمت راست، آدرس دهی تصادفی یک به یک را نشان می دهد.

2.4 معماری پردازنده های گرافیکی GeForce سری 8

پردازنده های گرافیکی سری 8800 که موسوم به G80 هستند در آبان ماه سال 1385 همراه با ارائه کارت های گرافیکی GeForce 8800 GTX مطابق با معماری کودا ساخته و به بازار ارائه شدند. یکی از مهم ترین تغییراتی که در معماری G80 اتفاق افتاده بود این بود که در این پردازنده ی گرافیکی برای نخستین بار واحد پردازش های رأسی و پیکسلی با یکدیگر ادغام شدند و یک واحد یکتا به وجود آمد که قادر بود همه ی پردازش های رأسی¹، پیکسلی²، هندسی و پردازش های موازی را باهم انجام دهد. دلایلی که شرکت NVIDIA را بدین سمت پیش برد در ادامه توضیح داده شده است.

- **توزیع بار ناهمگون:** یکی از مهم ترین علل ادغام این دو واحد وجود توزیع بار ناهمگون است. در پردازش تصاویر، زمانی که با زوایای بزرگ روبه رو هستیم، واحد پردازنده ی رئوس در بسیاری از مواقع بیکار است و بار زیادی بر دوش واحد پردازش پیکسلی قرار دارد. در عوض در زمان هایی که مشغول پردازش زاویه های کوچک هستیم، بار زیادی بر دوش واحد پردازش رأس قرار دارد و واحد پردازش پیکسلی بیکار است. در صورتی که برای هر کدام از این واحدها یک سخت افزار مجزا تعریف کنیم، نمی توانیم هر دو واحد را زیر بار کامل ببریم پس طراحی هرگز بهینه نخواهد بود.
- **کاهش هزینه:** دو واحد تلفیق شده می توانند از بسیاری از بخش ها به صورت مشترک استفاده کنند. این استفاده مشترک مساحت تراشه و هزینه های طراحی را کاهش می دهد.
- **طراحی یکپارچه:** با تلفیق این دو واحد یک گروه طراحی می تواند به راحتی بر روی جنبه های مختلف طراحی تمرکز نماید و طرحی بهتر و یکپارچه تر ارائه دهد.

¹ Vertex process

² Pixel process

The diagram illustrates the internal architecture of a GPU. At the top, the **Host CPU** connects to a **Bridge**, which in turn connects to **System memory**. The **Bridge** also connects to the **Host interface**. Below the **Host interface**, the architecture is divided into three main functional areas:

- Input assembler** leads to **Vertex work distribution**.
- Viewport/clip/setup/faster/z/cull** leads to **Pixel work distribution**.
- Compute work distribution** handles general-purpose computing tasks.

These distribution units feed into eight parallel processing units, each labeled **TPC** (Texture Processing Cluster). These TPCs are connected to a central **Interconnection network**. Below this network, the architecture includes six pairs of **ROP** (Render Output Processor) and **L2** (Level 2 cache) blocks, which are connected to **DRAM** (Dynamic Random Access Memory) blocks.

1- آرایه پردازنده‌های جریان (SPA')

2- رابطہ میزبان^۳

⁴Device

دارد. این واحد باید دستورات CPU را دریافت کرده، اطلاعات موردنیاز را از حافظه فراخواند و سازگار بودن دستورات مختلف را بررسی کند.

3- واحد Input Assembler

وظیفه ی این واحد واکنشی اطلاعات مربوط به تصاویر از حافظه است. این واحد قادر است در هر پالس ساعت یک مشخصه ی هندسی و هشت مشخصه ی رأسی را از حافظه برای پردازش واکنشی کند.

4- واحد Viewport/Clip/Setup/Raster/Zcull

این واحد وظیفه دارد تا یک سری پیش پردازش مقدماتی بر روی تصاویر انجام دهد. مهم ترین مسئولیت این قسمت تبدیل کردن اطلاعات برداری تصاویر به اطلاعات پیکسلی برای پردازش تصویر است.

5- واحدهای توزیع کار

واحدهای توزیع کار رأسی^۱، پیکسلی^۲ و پردازشی^۳ مسئولیت توزیع پردازش های مختلف را در بین واحدهای TPC بر عهده دارند. وظیفه ی این واحدها توزیع بار یکنواخت بین TPC هاست.

6- واحد ROC^۴

هر واحد ROC به یک قسمت مشخص از حافظه دسترسی دارد. TPC اطلاعات را از طریق شبکه ی اتصالات داخلی^۵ برای ROC ارسال می کند. وظیفه ی ROC کنترل عمق رنگ، تست رنگ، مخلوط کردن و به روزرسانی رنگ هاست. واضح است که تمامی این واحدها به صورت موازی باهم کار می کنند.

7- DRAM

در پردازنده های GeForce سری 8 حافظه ی RAM به صورت 8 بانک مستقل در نظر گرفته شده است. این سبک طراحی RAM کمک می کند تا پهنای باند بیشتری در اختیار پردازنده ها قرار گرفته و عملیات ذخیره و بازیابی به صورت موازی انجام شود. در معماری سری 8800 انتقال اطلاعات از RAM به واحد SPA از طریق کش L2 صورت می گیرد. تعداد پایه های باس داده در RAM برابر با 384 پایه است که به شش گروه که هر کدام دارای 64 پایه ی داده هستند تقسیم می گردد. شایان ذکر است که حافظه ی محلی^۶ و حافظه ی سرا سری^۷ بر روی این

¹Vertex Work Distribution

²Pixel Work Distribution

³Compute Work Distribution

⁴Raster Operation Processor

⁵Interconnection Network

⁶Local Memory

⁷Global Memory

RAM قرار دارند. حافظه ی محلی یک حافظه ی انحصاری است که برای هر ترد به صورت مستقل قابل دسترسی است. به عبارت دیگر تردهای مختلف نمی توانند به حافظه ی محلی همدیگر دسترسی پیدا کنند. حافظه ی عمومی هم حافظه ای است که برای همه ی تردها قابل دسترسی است. تردهای مختلف می توانند اطلاعاتشان را در حافظه ی عمومی باهم تبادل کنند.

8- واحد TPC

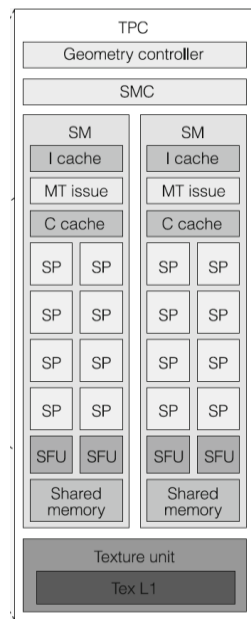
همان طور که در شکل 0-6 نشان داده شده است، هر واحد TPC دارای یک کنترل کننده ی اطلاعات هندسی¹، یک واحد کنترل کننده ی SM²، یک واحد بافت³ و دو SM⁴ است. واضح است که تعداد المان هایی که در داخل هر SM قرار دارد می تواند برحسب نیاز تغییر کند. واحد کنترل کننده ی اطلاعات هندسی موظف است خط لوله ی منطقی را که برنامه ی پردازش تصویر بر اساس آن نوشته شده است بر روی سخت افزار SM ها بنگارد. واحد SMC موظف است که SM ها مختلف را مدیریت کرده و بین آنها برای دسترسی به منابع مشترک از قبیل واحد بافت، دسترسی به حافظه و ورودی ها و خروجی ها داوری نماید. واحد SMC مسئول توزیع عادلانه بار محاسباتی بین SM ها نیز است. در یک TPC چند SM نیز وجود دارد که این SM ها وظیفه ی اجرای پردازش های رأسی، پیکسلی، هندسی و پردازش های موازی عمومی را بر عهده دارند. آخرین واحدی که در یک TPC قرار دارد واحد بافت است. واحد بافت قادر است در هر سیکل یک دستورالعمل را اجرا نماید. ورودی واحد بافت مختصات پیکسل هایی است که قرار است بر روی آنها عمل پردازش صورت گیرد و خروجی آن اطلاعات فیلتر شده است. واحد بافت به صورت کاملاً مستقل و موازی نسبت به SM عمل می کند. هر واحد پردازش بافت دارای چهار واحد تولید کننده ی آدرس و هشت واحد فیلتر کننده است که همه ی آنها قادرند به صورت موازی کار کنند. واحد بافت کاملاً به صورت خط لوله ای پیاده شده است و دارای یک حافظه ی پنهان است که پیاده کردن فیلترهای محلی را تسریع می کند.

¹Geometry Controller

²SM Controller (SMC)

³Texture Unit

⁴Streaming Multiprocessor



شکل 6-0: TPC

مهم ترین بخش TPC واحد SM^۱ است. این واحد دارای یک کش دستورالعمل^۱، یک کش فقط خواندنی ثوابت^۲، یک واحد واكشی و اجرای دستورالعمل^۳، 16 کیلوبایت حافظه ی مشترک خواندنی و نوشتنی، هشت پردازنده جریان^۴، و دو واحد پردازش های ویژه^۵ است. هر کدام از SP ها دارای یک واحد MAD^۶ است که عملیات ضرب و جمع متوالی به فرم $A = B \times C + D$ را در یک پالس ساعت انجام خواهد داد. SP ها یا پردازنده های جریان به تدریج با نام هسته های کودا شناخته شدند. پس از این ما نیز برای خطاب کردن پردازنده های جریان (SP) لفظ هسته های کودا را به کار خواهیم برد.

در هر SM واحد SFU دو وظیفه ی عمده بر عهده دارد. وظیفه ی نخست اجرای عملیات ریاضی سطح بالا مانند عملیات محاسبه ی توابع مثلثاتی و جذر بوده و وظیفه ی دوم الحاق پیکسل ها برای نمایش است. آخرین قسمت یک SM حافظه ی مشترک است. حافظه ی مشترک به برنامه هایی که بر روی هسته های کودای موجود در یک SM اجرا می شوند این امکان را می دهد که بدون نیاز به دسترسی به حافظه ی خارجی به تبادل اطلاعات بپردازند. بدین سان ترافیک دسترسی به حافظه ی خارجی کاهش یافته و سرعت نهایی اجرای برنامه افزایش خواهد یافت. در GeForce 8800 فرکانس پالس ساعت هسته های کودا و واحدهای SFU برابر با 1.5 GHz است و قدرت محاسباتی هر SM در بیشترین حالت برابر با 36 GFLOPS است. می توان گفت کوچک ترین واحد پردازنده ای

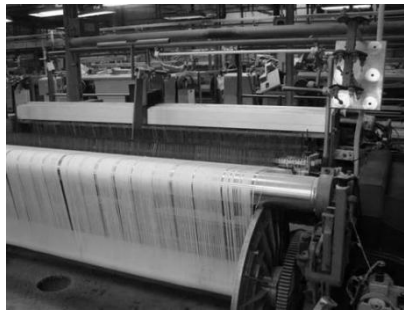
^۱Instruction Cache (I cache)^۲Constant Cache (C cache)^۳Multithread Instruction Fetch And Issue (MT issue)^۴Streaming Processor^۵Special Function Unit^۶Multiply ADD

که در یک پردازنده ی گرافیکی قرار دارد هسته های کودا است. تردهایی که برای اجرا به یک SM داده می شود بین هسته های کودا تقسیم شده و هر هسته ی کودا یک سری از آن ها را پردازش می کند به همین دلیل به این هسته های کودا پردازنده ی ترد هم گفته می شود. یک هسته ی کودا قادر است عملیات ریاضی ممیز شناور شامل جمع، ضرب و MAD انجام دهد. این عملیات با استاندارد IEEE 754 برای اعداد اعشاری ممیز شناور با دقت ساده منطبق بوده و قادر است مقادیر خاص از جمله NaN¹ یا مقادیر بی نهایت را تشخیص دهد. این هسته همچنین قادر است عملیات ریاضی صحیح متنوعی شامل مقایسه و تبدیل مقادیر انجام دهد. شایان ذکر است که این هسته ها به صورت کاملاً خط لوله ای طراحی شده اند.

یکی از معایب پردازنده های سری 8800 این است که در آن ها از اعداد زیر نرمال² پشتیبانی نمی شود. به عبارت دیگر در این پردازنده ها اعداد زیر نرمال با مقدار صفر تقریب زده شده و دور ریخته می شوند که این عمل دقت نهایی محاسبات را پایین خواهد آورد.

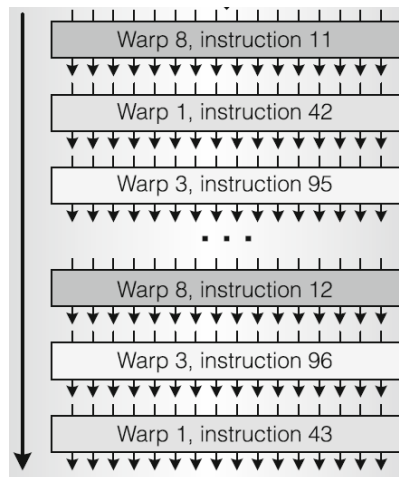
در کودا، کرنل برنامه ی ساده ای است که مشخص می کند یک ترد چگونه نتایج مورد نظر را محاسبه کند. برای اینکه پاسخ یک مسئله داده شود باید تردهای بسیار زیادی به صورت همزمان اجرا شوند. حال مسئله ای که ایجاد می شود این است که چگونه می توان اجرای موازی تردهای بسیار زیادی را که ممکن است مربوط به برنامه های مختلفی باشد مدیریت کرد. از آنجایی که هر برنامه ای که قرار باشد این تردها را مدیریت کند خودش زمان گیر است شرکت NVIDIA برنامه ی مدیریت اجرای این تردها را به صورت سخت افزاری پیاده کرده است این برنامه قادر است اجرای 768 ترد همزمان را با سربرار³ زمانی صفر ثانیه مدیریت کند. برای اینکه از برنامه های مستقل پردازش رئوس، گرافیک و برنامه های پردازش موازی پشتیبانی شود، هر SM قادر است مستقل از سایر SM ها برنامه ی مشخصی را اجرا کند مشکلی که ممکن است ایجاد شود این است که اجرای قسمت های مختلف برنامه بر عهده ی SM های مختلفی گذاشته شود و جایی که نیاز به تبادل داده است کار برخی از SM ها به پایان نرسیده باشد. برای حل چنین مشکلی دستورالعملی برای همزمان کردن SM در مواقع نیاز در نظر گرفته شده است. در پردازنده های گرافیکی NVIDIA از یک معماری جدید به نام SIMT⁴ استفاده می شود. واحد مدیریت اجرای تردهای SM، تردها را در دسته های 32 تایی که به آن وارپ⁵ گفته می شود اجرا می کند. واژه ی وارپ هم مثل ترد از صنعت نساجی به عاریه گرفته شده است و به معنی تار پارچه است. تار پارچه مجموعه ای از یک سری نخ موازی با هم است و اگر به شکل 7-0 توجه کنید متوجه خواهید شد که علت این نام گذاری چیست.

¹Not A Number²Sub normal³Overhead⁴Single Instruction Multiple Thread⁵Warp



شکل 7-0: نمونه ی تار پارچه

هر SM قادر است 24 وارپ را به صورت همزمان اجرا کند. تمامی تردهایی که در یک وارپ قرار دارند مشابه هم هستند و باید از یک آدرس برنامه شروع به اجرا کنند اما آنها آزادند که در طول اجرای برنامه به قسمت های مختلفی جهش کرده و به صورت مستقل اجرا شوند؛ هرچند که اگر چنین اتفاقی رخ دهد اجرای آنها به صورت موازی نخواهد بود. به این سبک اجرای موازی تردها SIMT گفته می شود. در هر SM واحد اجرای SIMT بررسی می کند و یکی از 24 وارپی را که برای اجرا آماده است انتخاب می کند. اولین انتخاب بر اساس نوع وارپ، نوع دستورالعمل و رعایت منصفانه بودن شانس اجرا برای همه ی وارپ ها صورت می گیرد. در معماری GeForce سری 8 هر وارپ به صورت دو سری 16 تردی اجرا می گردد که در مجموع 4 سیکل طول می کشد. نمونه ی اجرای SIMT در شکل 8-0 نشان داده شده است.



شکل 8-0: اجرای SIMT

پس از اینکه یک وارپ برای اجرا انتخاب شد SM، تردهای آن را بر روی هسته های کودا به صورت مستقل از هم اجرا می کند. بهترین بازده، زمانی به دست می آید که برنامه ای که قرار است اجرا گردد دارای انشعاب وابسته به داده نباشد. در صورتی که چنین چیزی وجود داشته باشد اصطلاحاً می گویند واگرایی اتفاق افتاده است. هر زمان که در

یک SM واگرایی اتفاق بیفتد تردهای آن به صورت متوالی اجرا خواهند شد. اگر بخواهیم دو روش SIMT و SIMD را باهم مقایسه کنیم باید گفت که روش اول کنترل ترد را بر عهده دارد ولی روش دوم کنترل برداری از داده ها را بر عهده می گیرد. تفاوت مهمی که وجود دارد این است که SIMT نه تنها اجازه ی نگارش برنامه ی موازی در سطح داده را می دهد بلکه اجازه ی اجرای تردهای مستقل را نیز فراهم می کند.

کامپایلرهای سطح بالا برنامه های گرافیکی و برنامه های پردازش موازی را مستقیماً به برنامه ی باینری قابل اجرا توسط پردازنده ی گرافیکی تبدیل نمی کنند بلکه آن را به یک سری دستورالعمل میانی مانند دستورالعمل های PTX تبدیل کرده و سپس این دستورالعمل ها پس از بهینه سازی به دستورالعمل های قابل اجرا توسط پردازنده ی گرافیکی تبدیل می شوند. بدین ترتیب شرکت می تواند نسل های متفاوتی از پردازنده های گوناگون تولید کند ولی زبان برنامه نویسی را تغییر ندهد. لازم به ذکر است که تغییر مداوم زبان برنامه نویسی، کاربران را نسبت به استفاده از سخت افزار دلسرد می کند زیرا کاربر هر روز مجبور است زبان جدیدی را بیاموزد. فایده ی دیگر این سبک کامپایلر برنامه این است که در آن می توان برنامه را قبل از تبدیل نهایی به کدهای باینری بررسی کرد، وابستگی های داده ها را بررسی کرد، محل های وقوع واگرایی را تشخیص داد، کدهای مرده (قسمت هایی که هرگز اجرا نخواهند شد) را کشف کرد و برنامه را تا حد امکان برای اجرا بر روی پردازنده ی گرافیکی بهینه کرد.

هر پردازنده ای متناسب با کاربردی که دارد دارای یک مجموعه دستورالعمل خاص است که در راستای هدف آن پردازنده بهینه شده اند. در ادامه مجموعه ی دستورالعمل هایی که توسط پردازنده های گرافیکی سری 8800 شرکت NVIDIA پشتیبانی می شوند آورده شده است.

• مجموعه دستورالعمل های اعشاری و صحیح

○ این دستورها شامل جمع، ضرب، ضرب-جمع یا همان MAD، یافتن کمینه، بیشینه، مقایسه و دستورهای تبدیل بین اعداد صحیح و اعشاری هستند. همچنین یک سری دستورات دیگر هم هست که توسط واحد SFU اجرا می شود. این دستورات عبارتند از محاسبه ی سینوس، کسینوس، مقدار نمایی، لگاریتم باینری، معکوس کردن عدد و محاسبه ی جذر.

○ این عملیات قادرند پرچم های^۱ صفر، منفی، کری^۲ و سرریز^۳ را که به صورت اختصاصی در هر ترد قرار دارند برای استفاده های بعدی ست کنند.

^۱ Flag

^۲ Carry

^۳ Overflow

- عملیات بیتی

- عملیات بیتی شامل شیفت به سمت راست، شیفت به سمت چپ، عملیات منطقی و جابه جایی هستند.

- عملیات کنترل اجرا

- این عملیات شامل انشعاب، فراخوانی، بازگشت، دستورالعمل trap و دستور مربوط به سنکرون کردن تردها است.

- دستورالعمل های دسترسی به حافظه

- برای دسترسی به حافظه دو نوع دستورالعمل در نظر گرفته شده است. نوع نخست دستورالعمل های بافت است که با کمک خود واحد بافت می تواند به اطلاعات حافظه دسترسی پیدا کند. نوع دیگر دستورالعمل های ذخیره و بازیابی معمولی است. این دستورالعمل ها برای هماهنگ شدن کامل سخت افزار با زبان C به آن اضافه شده است. این دستورها توسط SM اجرا می گردند.

تا اینجا با معماری پردازنده های گرافیکی سری 8 و معماری کودا آشنا شدیم. اکنون زمان مناسبی است که بین این دو معماری ارتباط برقرار کنیم و ببینیم که مدلی که کودا برای برنامه نویسی ارائه می کند چگونه بر روی سخت افزار پردازنده های گرافیکی نگاشته می شود.

گفته شد که برای اجرای یک برنامه ی موازی نیاز است تعداد بسیار زیادی ترد اجرا شوند و هر کدام از آن ها کار خاصی را انجام دهند تا نهایتاً پاسخ مسئله به دست آید. همچنین گفته شد که این تردها ممکن است در SM های گوناگون باشند و هر SM این تردها را در قالب هایی به نام وارپ اجرا می کند. فرض کنید فقط یک برنامه بر روی پردازنده ی گرافیکی در حال اجرا است؛ برای اجرای این برنامه یک سری ترد در حال اجرا شدن هستند به این تردها که همه مربوط به اجرای یک برنامه بوده و برای یافتن پاسخ ممکن است به مبادله ی اطلاعات نیاز داشته باشند¹ CTA گفته می شود. تعداد تردهای درون یک CTA می تواند از 1 تا 512 عدد باشد. هر کدام از این تردها دارای یک شماره ی منحصر به فرد در درون CTA هستند و برنامه نویس می تواند تردها را درون CTA به صورت یک بعدی، دوبعدی و یا سه بعدی بچیند. تردهایی که درون یک CTA هستند قادرند داده هایشان را از طریق حافظه ی سراسری و یا محلی باهم مبادله کنند و یا در نقطه های خاصی از اجرا باهم سنکرون گردند. هر SM قادر

¹Cooperative Thread Array

است نهایتاً تا هشت CTA را به صورت هم زمان اجرا کند. این تعداد وابسته است به منابعی که هر CTA نیاز دارد. واحد SMC در صورتی که ببیند منابع خالی وجود دارد یک CTA تعریف کرده و تردهای جدیدی را اجرا می کند. واضح است که SM، هر CTA را در قالب وارپ هایی که هر کدام دارای 32 ترد است اجرا خواهد کرد. این طراحی خارق العاده به اینجا ختم نمی شود؛ در پردازنده های سری 8 مجموعه ای از CTA ها قرار دارد که به آن گرید CTA گفته می شود. در این گرید CTA های زیادی قرار دارد که واحد توزیع کار این CTA ها را با توجه به اینکه کدام SM دارای ظرفیت خالی است بین آن ها توزیع می کند. برای اینکه بتوان فارغ از اندازه پردازنده ی گرافیکی به راحتی برنامه نویسی کرد، CTA ها کاملاً مستقل از هم اجرا می شوند و تردهایی که در داخل CTA های مختلف قرار دارند امکان تبادل اطلاعات را با هم نخواهند داشت. بدین سان ممکن است تمام CTA های یک برنامه ی خاص بر روی یک پردازنده ی گرافیکی به صورت موازی با هم اجرا شوند و بر روی پردازنده ی گرافیکی دیگر هم به صورت متوالی اجرا شوند در این حالت تنها تفاوتی که حاصل می شود در زمان اجرای برنامه است و اشتباهی در خروجی پدید نخواهد آمد. پس سلسله مراتب اجرا در پردازنده ی گرافیکی دارای سه قسمت زیر است:

- ترد: یک بخش از پاسخ را بیان می کند و با TID مشخص می گردد.
- CTA: پاسخ بلاک هایی از برنامه را مشخص می کند و با CTA ID مشخص می شود.
- Grid: پاسخ تعداد زیادی CTA را مشخص می کند.

با این تفاسیر تطبیق مدل برنامه نویسی کودا با معماری سخت افزاری کار دشواری نیست. هر ترد در مدل برنامه نویسی کودا با یک ترد در سخت افزار معادل است. هر بلاک کودا معادل CTA و گرید کودا با گریدی که در این قسمت معرفی شد برابر است. پس بلاک هایی که در کودا تعریف می شود در قالب CTA در SM اجرا خواهد شد.

2.5 معماری پردازنده های گرافیکی GeForce سری 200

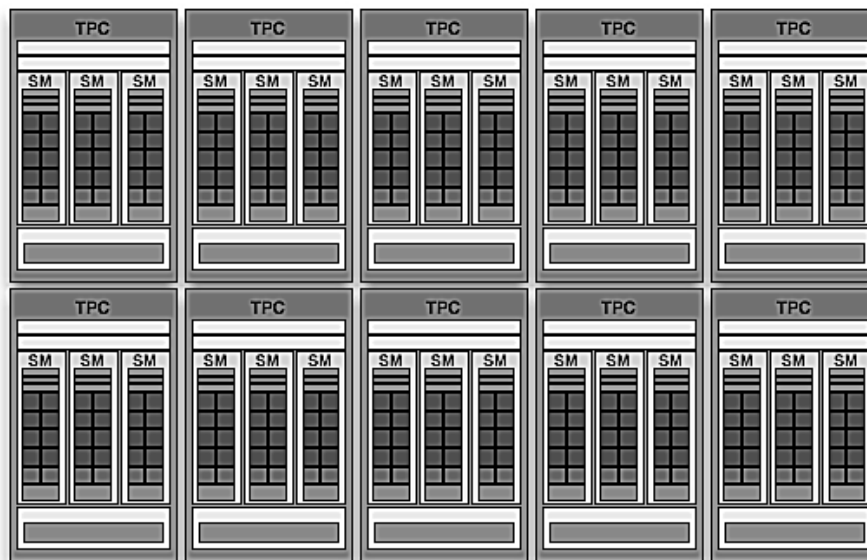
پردازنده های گرافیکی سری 200 دومین نسل پردازنده های گرافیکی شرکت NVIDIA با معماری یکپارچه واحد سایه زنی هستند. شرکت NVIDIA پردازنده های گرافیکی سری 200 را با هدف دستیابی به اهداف زیر تولید کرد.

- بالا بردن کارایی پردازنده های گرافیکی نسل قبل یعنی GeForce سری 8
- بهینه سازی معماری برای بازی های رایانه ای جدید

- بهینه سازی معماری برای کاهش مصرف انرژی و فضای مصرفی

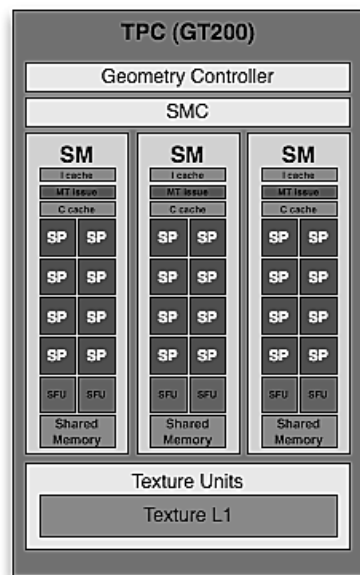
- ساخت یک پردازنده ی گرافیکی قدرتمند برای اجرای پردازش موازی با مقیاس بزرگ تر

پردازنده ی گرافیکی GeForce GTX 200 اولین پردازنده ای است که از نسل دوم واحد سایه زنی یکپارچه شرکت NVIDIA استفاده می کند. این پردازنده دارای سرعتی به صورت متوسط 1.5 برابر سرعت پردازنده های نسل های پیشین است. معماری این پردازنده بر اساس معماری پردازنده های گرافیکی GeForce سری هشت است با این تفاوت که معماری این پردازنده در بسیاری از بخش ها بهینه شده است. در پردازنده های سری 200 تعداد TPC ها از 8 به 10 عدد در هر پردازنده افزایش پیدا کرده است و تعداد SM های موجود در هر TPC نیز به سه عدد افزایش پیدا کرده است. پس با یک حساب سرانگشتی می توان متوجه شد که تعداد هسته های کودا در این نسل از پردازنده ها به عدد 240 رسیده است. در تمامی پردازنده های گرافیکی سری 200 نزدیک به 80 درصد از ترانزیستورها برای مصرف پردازشی استفاده شده است در حالی که در CPU ها به صورت متوسط نزدیک به 20 درصد از فضا برای پردازش اختصاص داده شده و حجم زیادی از فضا برای حافظه ی کش اختصاص پیدا کرده است. از این مقایسه می توان نتیجه گرفت که در طراحی CPU ها تمرکز بسیار بالایی بر روی کاهش تأخیر وجود دارد؛ انتظار می رود سرعت پاسخ یک CPU بسیار بالا باشد. این در حالی است که در مورد پردازنده های گرافیکی تمرکز بر روی بالا بردن قدرت پردازش موازی بر روی حجم های بسیار زیاد داده های ورودی است.



شکل 9-0: معماری GeForce GT 200

در شکل 9-0 معماری GeForce GT 200 نشان داده شده است. در بالاترین قسمت واحد برنامه ریز ترد¹ وجود دارد که اجرای تردهای مختلف بین TPC ها را مدیریت می کند. در قسمت پایین هم واحدهای بافت و واحد دستورالعمل های اتمیک² وجود دارد که نقش آن ها مانند نقش همین قسمت ها در پردازنده های نسل قبل است. همان طور که شکل نشان می دهد، این سری پردازنده ها دارای ده واحد TPC است که اجزای این واحد، با جزئیات در شکل 10-0 نشان داده شده است. همان طور که در شکل نشان داده شده است اجزای واحد TPC مانند اجرای واحد TPC در نسل قبلی پردازنده ها است، با این تفاوت که در این نسل تعداد SM بیشتر شده است. پردازنده های سری 200 هم مانند نسل قبلی نسبت به تأخیر در فراهم شدن داده مقاوم هستند؛ یعنی در صورتی که یک ترد برای دسترسی به داده معطل بماند، سخت افزار قادر است بدون تحمیل کردن هیچ گونه سرباری مجموعه تردهای دیگری را اجرا کند.



شکل 10-0: اجزای واحد TPC در پردازنده های گرافیکی سری 200

یکی از ایرادهای پردازنده های سری های قبل، این بود که در اجرای برنامه های بلند به مشکل پر شدن رجیسترها برخورد می کردند. در این نسل از پردازنده های گرافیکی اندازه فایل رجیسترها دو برابر شده است که نوشتن برنامه هایی طویل با کارایی بالا را تسهیل خواهد کرد. یکی دیگر از مزایای معماری GT 200 این است که واحد محاسبات اعشاری آن محاسبات ضرب جمع را بدون گرد کردن میانی محاسبه می کند. در پردازنده های قبلی محاسبات MAD در دو مرحله گرد می شد. مرحله ی نخست پس از محاسبه ی حاصل ضرب بود و مرحله ی دوم پس از محاسبه ی حاصل جمع. این گرد کردن دومرحله ای باعث افت دقت می شد. در پردازنده های سری 200

¹Thread Scheduler

² Atomic

عملیات MAD روی اعداد اعشاری با دقت مضاعف، حاصل را پس از عملیات ضرب گرد نکرده و آن را با حفظ تمامی رقم های اعشاری با مقدار ثابت نهایی جمع خواهد کرد. به این روش محاسبه که نسبت به روش قبلی از دقت بالاتری برخوردار است Fused Multiply ADD گفته می شود.

پردازنده های GeForce سری 200 در بسیاری از جزئیات از جمله واحدهای بافت، محاسبات هندسی، واحدهای ROC، مدیریت حافظه و مدیریت مصرف انرژی نسبت به نسل های پیش از خود برتری داشتند که بیان به تفصیل این جزئیات از موضوع بحث این کتاب خارج است. پردازنده های گرافیکی GeForce سری 200 با توانایی های بالایی که داشتند پردازش موازی همه منظوره با کمک پردازنده های گرافیکی را بیش از پیش فراگیر کردند. پس از پردازنده های این نسل، شرکت NVIDIA معماری نسل سوم پردازنده های گرافیکی با واحد سایه زنی یکپارچه را تحت عنوان فرمی^۱ ارائه کرد.

2.6 معماری Fermi

زمانی که شرکت NVIDIA شروع به توسعه معماری فرمی کرد، پردازش موازی با کمک پردازنده های گرافیکی دیگر امری نوپا نبود. در آن زمان برنامه های بسیاری برای دو نسل قبلی از پردازنده های گرافیکی توسعه پیدا کرده بود و نظرات بسیاری از سوی کاربران برای شرکت ارسال شده بود. بدین ترتیب NVIDIA توانست با بهره گیری از تجربیات حاصل از پردازنده های گرافیکی دو نسل قبل یک معماری قدرتمند ارائه دهد. در این معماری کاستی های نسل های قبل برطرف شد و استفاده از پردازنده ی گرافیکی برای توسعه ی برنامه های موازی آسان تر و بهینه تر گردید. مهم ترین ویژگی هایی که از دیدگاه پردازش موازی در معماری فرمی رعایت شده است در ادامه ذکر شده است.

- بهینه سازی کارایی واحد محاسبات اعشاری با دقت مضاعف^۲

○ هرچند واحد محاسبات اعشاری با دقت معمولی در پردازنده های گرافیکی قبلی دارای کارایی قابل قبولی بود اما برخی از مسائل مربوط به پردازش های موازی نیازمند کارایی بالا در کار با اعداد اعشاری با دقت مضاعف هستند.

- استقرار واحد تشخیص و تصحیح خطای ECC^۳

^۱Fermi

^۲Double Precision Floating Point

^۳Error Correcting Code

○ با کمک واحد ECC می توان تعداد بسیار زیادی GPU را در یک مرکز پردازش موازی در کنار هم قرارداد و به راحتی و بدون نگرانی از بروز خطا بر روی آن ها، محاسبات حساب به خطا از جمله محاسبات مربوط به امور مالی را انجام داد.

• کش

○ برخی از الگوریتم های پردازش موازی نمی توانند از حافظه ی مشترکی که بر روی پردازنده ی گرافیکی در نظر گرفته شده است استفاده کنند به همین دلیل امکان استفاده از کش واقعی نیز در معماری فرمی در نظر گرفته شده است.

• افزایش حجم حافظه مشترک^۱

○ بر اساس تجربه های به دست آمده از نسل های قبلی پردازنده های گرافیکی، مشخص شده که برخی از امور پردازش موازی نیاز به بیش از 16 کیلوبایت حافظه مشترک دارند تا بتوانند به سرعت بالاتری برسند.

• تغییر محتوای سریع

○ پردازنده هایی که در داخل پردازنده ی گرافیکی قرار دارند هم مسئول رسیدگی به امور گرافیکی هستند و هم باید پردازش های مربوط به پردازش موازی را بر عهده بگیرند. بدین ترتیب لازم است که بارها کار خود را تغییر دهند. این تغییر کار^۲ فرآیندی زمان بر است که پردازنده را ملزم می سازد تا تمامی محتوای داخلی خود را با محتوای جدیدی تغییر دهد. در معماری فرمی امکان تغییر این محتوا با سرعت بالاتری فراهم شده است.

• دستورهای اتمیک سریع تر

○ بر اساس تجربه های قبلی نیاز به دستورهای اتمیک سریع تری برای پیاده سازی الگوریتم های پردازش موازی بر روی پردازنده های گرافیکی وجود داشت. در معماری فرمی این دستورات نیز گنجانده شده است. دستورالعمل های اتمیک به دستورالعمل هایی گفته می شود که اجازه ی انجام عملیات بازایی، پردازش و ذخیره مجدد قسمت های مختلفی از یک ساختمان داده را به صورت هم زمان به تردهای مختلف می دهد. عملیاتی اتمیک نامیده می شوند که شامل

^۱Shared Memory

^۲Task Switching

دستورهای بازیابی، پردازش و ذخیره برای هر ترد به صورت مستقل و بدون ایجاد مزاحمت با سایر تردها هستند. به دلیل افزایش تعداد واحدهای اتمیک در فرمی و وجود کش L2 سرعت عملیات اتمیک در فرمی نسبت به GT200 بیش از بیست برابر بیشتر شده است.

این ویژگی های فرمی به علت وجود یک معماری منحصربه فرد با اجزای مناسب و متناسب است. در ادامه برخی از واحدهایی که در بهبود معماری فرمی نقش داشته اند توضیح داده خواهد شد.

• افزودن SM های نسل سوم

○ در نسل سوم SM ها تعداد هسته های کودا در هر SM به سی و دو هسته افزایش پیدا کرد که نسبت به معماری GT200 چهار برابر شده بود. همچنین توانایی محاسبات اعشاری با دقت مضاعف نسبت به پردازنده ی گرافیکی نسل قبل 8 برابر افزایش یافت. از جمله ی دیگر بهینه سازی هایی که در معماری فرمی پدید آمد اضافه شدن قسمت مدیریت وارپ مضاعف بود که به پردازنده این امکان را می داد که بتواند دستورالعمل ها را به صورت مستقل بین دو وارپ مختلف به صورت هم زمان توزیع و اجرا نماید. بالاخره اینکه در معماری فرمی 64 کیلوبایت حافظه ی رم اضافه شد که به صورت پویا قابل تقسیم به دو قسمت 16 و 48 کیلوبایتی بین کش L1 و حافظه ی مشترک بود.

• نسل دوم دستورالعمل های موازی اجرای ترد

○ در این نسل فضای آدرس را برای تطبیق پذیری کامل با C++ بهینه کرده اند و همچنین واحد محاسبات اعشاری با دقت ساده و مضاعف ایجاد شد که کاملاً با استاندارد IEEE 754 سال 2008 میلادی مطابقت دارد. همچنین در این نسل واحد پیش بینی انشعاب¹ ایجاد شده است که تأثیر به سزایی در افزایش کارایی پردازنده های گرافیکی دارد.

• بهینه سازی واحدهای حافظه

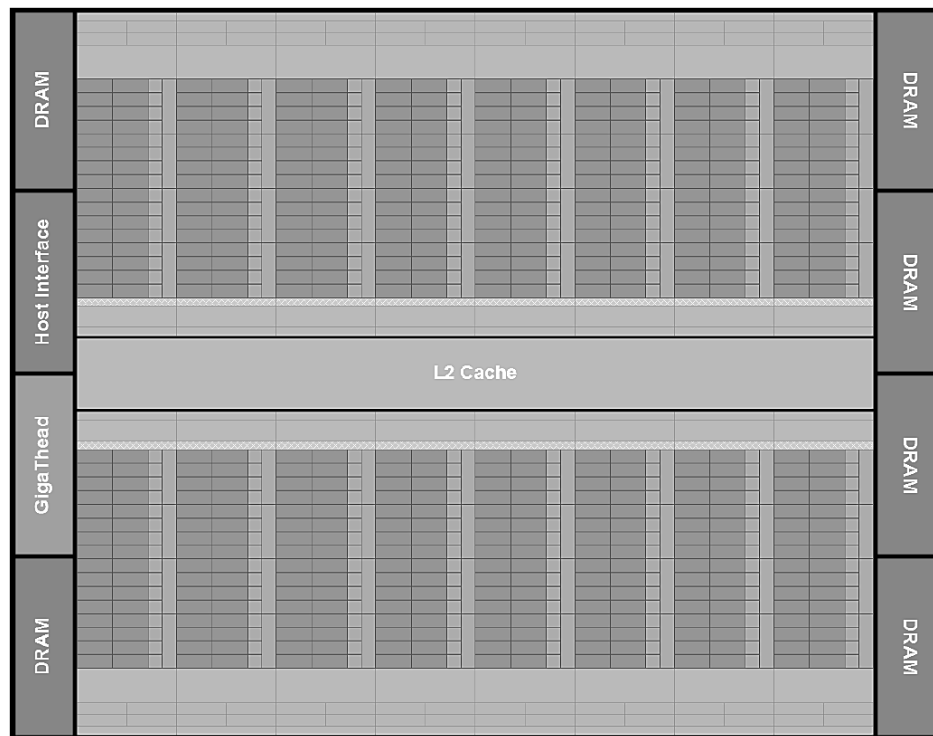
¹ Branch Prediction

○ در این نسل از حافظه ها سه تغییر مهم رخ داد. نخست افزوده شدن واحد تشخیص و کنترل خطا بود. تغییر دوم بهینه سازی ساختار کش و حافظه ی مشترک بود و سوم بهینه سازی دستورات اتمیک در این نسل بوده است.

• افزودن موتور GigaThread™

○ با کمک این موتور، فرمی قادر است با سرعت ده برابر معماری نسل قبلی خودش، تغییر محتوا را انجام دهد. همچنین به کمک این موتور توانایی اجرای همزمان کرنل ها و اجرای خارج از نوبت بلاک ها فراهم شد.

اولین پردازنده ی گرافیکی که بر اساس معماری فرمی طراحی شد دارای سه میلیارد ترانزیستور و 512 هسته ی کودا بود. هر هسته ی کودا قادر است یک دستورالعمل اعشاری یا صحیح را برای هر ترد در هر پالس ساعت اجرا نماید. این 512 هسته ی کودا در شانزده SM که هر کدام دارای 32 هسته ی کودا هستند قرار گرفته اند. این پردازنده دارای شش حافظه با پهنای باند داده ی 64 بیت است و مجموع این حافظه ها می تواند به 6 گیگابایت حافظه ی رم GDDR5 برسد. شکل 0-11 نمونه ای از یک پردازنده ی گرافیکی با معماری فرمی را نشان می دهد.



شکل 0-11: معماری فرمی

همان طور که شکل 0-11 نشان می دهد حافظه ی کش L2 که در وسط قرار گرفته است بین تمامی SMها مشترک است. هر کدام از مستطیل های عمودی نشان دهنده ی یک SM است. نوار پایین نشان دهنده ی واحد توزیع کننده و برنامه ریز تردهاست و مستطیل های کوچک واحدهای اجرایی هستند.

2.6.1 SMهای نسل سوم

نسل سوم SMها دارای ویژگی های بارزی است که آن را به صورت کامل از معماری نسل دوم جدا می کند. هر SM دارای 32 هسته ی کوداست. هر کدام از این هسته ها دارای یک واحد محاسبات منطقی و ریاضی¹ و یک واحد محاسبات اعشاری² است که هر دوی آنها دارای خط لوله کامل هستند. در نسل سوم SMها از استاندارد IEEE 754 سال 2008 به صورت کامل پشتیبانی می شود. این پردازنده ی گرافیکی دارای یک واحد FMA³ است که عملیات ضرب-جمع را فقط با یک بار گرد کردن پاسخ، محاسبه می کند که این باعث بهبود دقت محاسبات می گردد. در معماری GT 200 واحد محاسبات اعشاری دارای دقت بیست و چهار بیتی بود. در معماری فرمی این دقت به 32 بیت رسید که برای استفاده شدن در زبان های برنامه نویسی استاندارد مناسب بود. در این

¹ ALU² FPU³ Fused Multiply-Add

معماری دقت ALU هم به 64 بیت رسید و تنوع دستورالعمل ها افزایش پیدا کرد. این دستورالعمل ها شامل عملیات منطقی، شیفت، جابه جایی، مقایسه و قراردعی bit-reversed است. هر SM دارای 16 واحد ذخیره و بازیابی اطلاعات است که اجازه می دهد عملیات محاسبه ای آدرس به صورت همزمان برای 16 ترد در هر کلاک محاسبه شود. از طرف دیگر هر SM دارای 4 واحد SFU¹ بوده که قادرند دستورات سطح بالا را مانند سینوس، کسینوس، محاسبه ی معکوس و جذر را اجرا نمایند. هر SFU قادر است یک دستورالعمل را در هر کلاک اجرا نماید. پس دستورات یک وارپ را در هشت کلاک اجرا خواهد کرد. شایان ذکر است که تا زمانی که واحدهای محاسبات اعشاری مشغول است امکان استفاده از سایر واحدهای SM وجود دارد. محاسبات اعشاری با دقت مضاعف در بسیاری از مسائلی که نیاز به پردازش موازی دارند وجود دارد در معماری فرمی امکان اجرای 16 دستورالعمل اعشاری با دقت مضاعف در هر کلاک و به ازای هر SM وجود دارد. در معماری فرمی هر SM دو واحد مدیریت وارپ و دو واحد اجرای دستورالعمل دارد که پردازنده ی گرافیکی را توانمند می سازد تا بتواند به صورت همزمان دو وارپ را اجرا کند. این واحد دو وارپ را انتخاب می کند و از هر کدام از آن ها یک دستورالعمل را بر روی 16 هسته، یا 16 واحد ذخیره و بازیابی یا 4 واحد SFU اجرا می کند. واضح است که هر دو دستورالعملی نمی توانند به این ترتیب دوتایی باهم اجرا شوند. دو دستورالعمل صحیح، دو دستورالعمل اعشاری، یک دستورالعمل صحیح و یک دستورالعمل اعشاری، دستورهای ذخیره و بازیابی و دستورهای واحد SFU می توانند به صورت همزمان اجرا شوند.

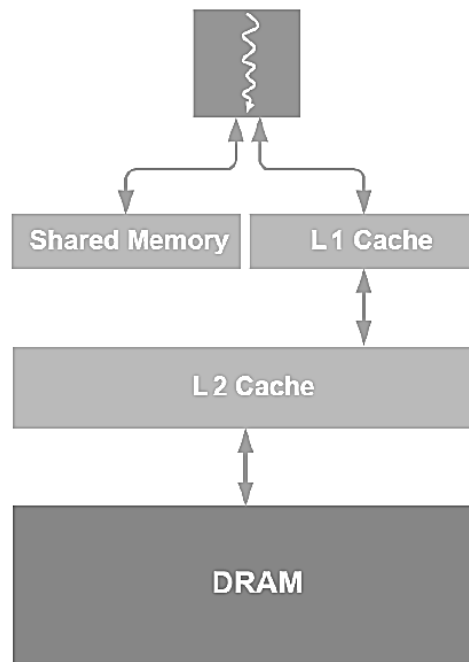
2.6.2 حافظه با قابلیت بازپیکربندی

حافظه ی مشترک نقش مهمی را در برنامه هایی که با پردازنده ی گرافیکی نوشته می شوند بازی می کند. این حافظه برای برقراری ارتباط بین تردهای مختلف در پردازنده ی گرافیکی مورد استفاده قرار می گیرد و از همه مهم تر نیاز به دسترسی به حافظه ی عمومی را کاهش می دهد. علیرغم فواید حافظه ی مشترک برخی از الگوریتم های پردازش موازی قادر نبودند از توانایی های این نوع حافظه به اندازه ی کافی بهره برداری کنند و نیاز به یک حافظه پنهان واقعی داشتند. این مشکل در پردازنده های فرمی حل شد. در فرمی هر SM دارای 64 کیلوبایت حافظه است که این حافظه بنا بر نیاز کاربر می تواند به دو قسمت 16 کیلوبایتی و 48 کیلوبایتی تقسیم شود. در صورتی که کاربر استفاده ی بیشتری از حافظه ی مشترک نسبت به کش بکند می تواند قسمت 48 کیلوبایتی را به حافظه ی مشترک اختصاص داده و قسمت 16 کیلوبایتی را به کش اختصاص دهد و در صورتی که کاربر به کش بیش از حافظه ی مشترک نیاز داشته باشد می تواند جای این دو را باهم عوض کند. شایان ذکر است که در معماری GeForce سری

¹ Special Function Unit

8 فضای که برای کش در نظر گرفته شده بود برابر با 16 کیلوبایت بوده است که در فرمی این فضا هم از لحاظ حجم و هم از لحاظ امکان بازپیکربندی بهبود یافته است.

مسئله دیگر این است که برخی از مسائل نیاز به داشتن یک حافظه ی کش دارند و برخی دیگر نیاز به داشتن حافظه ی مشترک و برخی هم نیاز دارند هر دو حافظه را داشته باشند. شرکت NVIDIA این مشکل را با ارائه دو سطح کش حل کرده است. در این روش یک کش L2 وجود دارد که بین تمامی MP ها مشترک است. این کش که دارای ظرفیتی معادل با 768 کیلوبایت است به SM ها این امکان را می دهد که به جای اینکه اطلاعات را در حافظه ی اصلی با یکدیگر مبادله کنند آن ها را در کش مبادله کنند و به سرعت بالاتری دست پیدا کنند. از طرف دیگر یک کش سطح یک نیز وجود دارد که سایز آن قابل تنظیم است. در حقیقت در این سطح یک حافظه با حجم 64 کیلوبایت وجود دارد که کاربر با توجه به نوع کاربرد می تواند آن را به دو قسمت 16 و 48 کیلوبایتی تقسیم کند. در صورتی که کاربر نیاز بیشتری به کار با کش احساس کند می تواند قسمت 48 کیلوبایتی را به کش اختصاص داده و قسمت 48 کیلوبایتی را برای حافظه ی مشترک در نظر بگیرد و بالعکس. لازم به ذکر است که این کش سطح 1 و حافظه ی مشترک در هر SM به صورت اختصاصی وجود دارد. پس تردهایی که در داخل یک SM قرار دارند نمی توانند به حافظه ی مشترکی که در داخل SM دیگری قرار دارد دسترسی پیدا کنند و یا اینکه از مزایای کش L1 ای که در داخل SM دیگر است بهره مند شوند. اگر بخواهیم به صورت خلاصه بیان کنیم، کش L1 و حافظه ی مشترک برای افزایش سرعت و تسهیل ارتباط داخل SM ها کمک می کنند در حالی که کش L2 در افزایش سرعت ارتباط بین SM ها کمک می کنند.



شکل 0-12: ساختار حافظه در فرمی

2.6.3 واحد تصحیح و تشخیص خطا

در صورتی که قرار باشد یک پردازنده در مسائل حساس مورد استفاده قرار بگیرد، باید حتماً صحت ثبت و انتقال اطلاعات روی باس های مدارهای آن بررسی شود. فرمی اولین پردازنده ی گرافیکی است که می تواند این بررسی را انجام دهد. فرمی که مجهز به واحد ECC¹ است که قادر است خطای هایی را که منجر به تغییر یک بیت شده اند را تصحیح و خطاهایی را که منجر به تغییر دو بیت شده اند را تشخیص دهد. پس در صورت بروز چنین خطاهایی این امکان فراهم خواهد شد که به جای اینکه برنامه با همان داده های غلط کار خودش را ادامه دهد متوقف شده و کار را از ابتدا آغاز کند. با توجه به اینکه چنین خطاهایی به علت وجود تشعشعات رادیویی صورت می گیرد باید توجه کرد که در کاربردهایی که در آن ها امکان وجود چنین خطاهایی وجود دارد از پردازنده های گرافیکی استفاده کرد که مجهز به فناوری ECC باشند. برای مثال برای استفاده از پردازنده ی گرافیکی در یک دستگاه MRI یا یک دیتاسنتر بزرگ باید از پردازنده های گرافیکی که دارای چنین قدرتی هستند بهره گرفت. در فرمی، رجیسترها، کش های L1 و L2 و حافظه ی DRAM دارای فناوری تصحیح و تشخیص خطا بوده و همچنین فرمی از استاندارد صنعتی برای کنترل خطا برای نقل و انتقال داده بین تراشه های مختلف نیز استفاده می کند. این بررسی های

¹Error Correcting Code

بروز خطا در قسمت های مختلف، پردازنده ی فرمی را برای استفاده شدن در سیستم های پردازش موازی که نسبت به خطا حساس هستند میسر می سازد.

2.6.4 تکنولوژی PTX 2.0¹

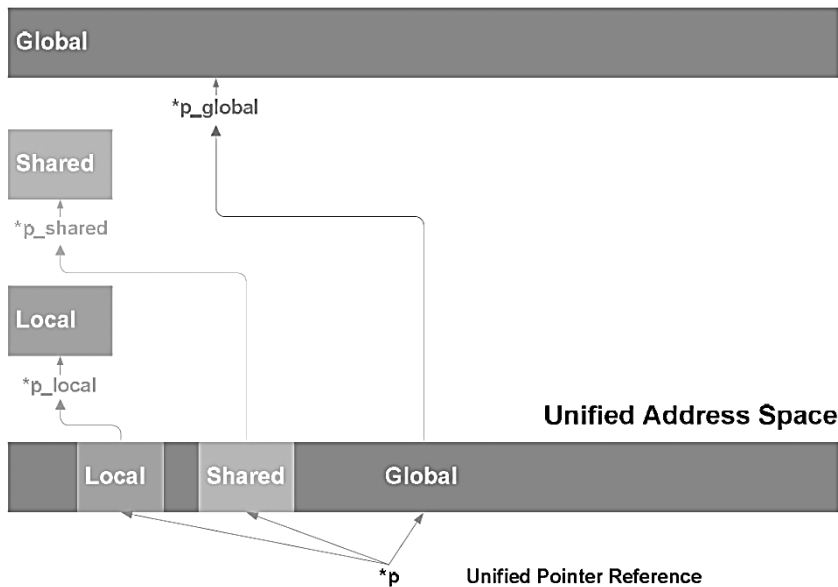
فرمی اولین پردازنده ای است که از تکنولوژی PTX 2.0 پشتیبانی می کند. PTX یک ماشین مجازی سطح پایین است که به صورت یک لایه واسطه بر روی پردازنده ی گرافیکی قرار می گیرد. در هنگام اجرای برنامه دستورالعمل هایی که برای PTX نوشته شده اند توسط درایور پردازنده ی گرافیکی به زبان ماشین تبدیل می شوند. علل اصلی که باعث شد تا شرکت NVIDIA از ماشین مجازی PTX استفاده کند عبارتند از:

- امکان ایجاد یک زبان برنامه نویسی برای خود PTX که به راحتی بتواند علیرغم تغییر نسل های مختلف پردازنده های گرافیکی بدون تغییر باقی بماند.
- امکان توسعه برنامه با زبان های برنامه نویسی مستقل از سخت افزار مانند C, C++, Fortran بر روی پردازنده ی گرافیکی
- فراهم کردن یک زبان برنامه نویسی که قادر باشد بر روی پردازنده ی گرافیکی که دارای تعداد هسته های گوناگونی هستند کار کرده و برنامه را متناسب با نیاز پردازنده ی گرافیکی توزیع کنند.

در PTX 2.0 برای نخستین بار فضای آدرس حافظه یکپارچه شد. پیش از این سه فضای مستقل آدرس حافظه وجود داشت که عبارتند از فضای آدرس حافظه ی انحصاری هر ترد، فضای آدرس حافظه ی مشترک داخل بلاک ها و فضای آدرس حافظه ی عمومی. در PTX 1.0 دستورالعمل های ذخیره و بازیابی هر کدام از این قسمت ها مستقل از سایر قسمت ها بود. به عبارت دیگر می بایست در زمان کامپایل کردن برنامه مشخص باشد که مقصد هر دستورالعمل کدام واحد از واحدهای حافظه است. این محدودیت باعث می شد که پیاده سازی اشاره گرهای C و C++ بر روی پردازنده های گرافیکی ممکن نباشد زیرا در زمان کامپایل برنامه مقصد اشاره گرها مشخص نیست. در فرمی فضای آدرس قسمت های مختلف با یکدیگر ادغام شد و نهایتاً یک فضای یکتا ایجاد شده است که از طریق آن می توان به تمام قسمت های حافظه دست پیدا کرد.

¹Parallel Thread eXecution (PTX) 2.0

Separate Address Spaces



شکل 0-13: فضای آدرس یکتا در معماری فرمی

در معماری فرمی قابلیت های کار با اعداد اعشاری هم بهبود پیدا کرد. در فرمی اعداد غیرمعمولی¹ به صورت پیش فرض توسط سخت افزار قابل شناسایی شد. همین طور چهار روشی که در سند IEEE 754 2008 برای گرد کردن اعداد در نظر گرفته شده است نیز در فرمی قابل پشتیبانی هستند. این چهار روش عبارت اند از: گرد کردن به نزدیک ترین عدد مجاور، گرد کردن به سمت صفر، گرد کردن به سمت مثبت بی نهایت و گرد کردن به سمت منفی بی نهایت. اعداد غیرمعمولی اعدادی هستند که بین صفر و کمترین عدد غیر صفر قرار دارند. در نسل های قبلی پردازنده های گرافیکی این اعداد صفر در نظر گرفته می شدند که باعث افت دقت می شد. در فرمی این اعداد توسط پردازنده های گرافیکی قابل پردازش هستند. شایان ذکر است که در CPU این اعداد توسط روتین های مدیریت خطا بررسی می شوند که بررسی آن ها زمان نسبتاً زیادی طول می کشد.

2.6.5 اجرای هم زمان کرنل ها

در معماری فرمی از اجرای هم زمان کرنل ها نیز پشتیبانی می شود یعنی کرنل های مختلف امکان اجرا شدن به صورت هم زمان را دارند. بدین ترتیب در هر زمان امکان استفاده از تمام قدرت پردازنده ی گرافیکی فراهم شده و سرعت نهایی اجرای برنامه بیشتر خواهد شد. این موضوع در شکل 0-14 نشان داده شده است.

¹Sub normal



شکل 2-14: اجرای همزمان کرنل ها در فرمی

2.6.6 کاستی ها فرمی

علیرغم کارایی بالای فرمی هنوز هم کاستی هایی وجود دارد که برای تبدیل کردن پردازنده ی گرافیکی به یک ماشین پردازش موازی کارا باید برطرف گردند. برخی از مهم ترین این کاستی ها عبارتند از:

1- حافظه ی نسبتاً کوچک پردازنده های گرافیکی

یکی از فرض هایی که در تمام سیستم های پردازش موازی فوق سریع مطرح می گردد این است که برنامه و تمامی اطلاعات مورد نیاز آن، همه با هم در حافظه ی پردازنده هستند. بر اساس قانون آمدال که در زمان توسعه ی mainframe بیان شد در یک سیستم پردازش فوق سریع لازم است 1 بایت حافظه و 1 بایت بر ثانیه ورودی و خروجی به ازای هر دستورالعمل بر ثانیه وجود داشته باشد پس در سامانه ای با توانایی پردازش 500 GFLOPS باید 500 GB حافظه وجود داشته باشد. هرچند ممکن است بسیاری تصور کنند که این حجم حافظه بیش از حد بزرگ است اما نباید این نکته را فراموش کرد که توسعه ی پردازنده های فوق سریع فقط برای حل مسائل دشوارتر است. مسائلی که به دلیل وجود ابعاد وسیع هنوز برای بشر قابل حل نیستند.

2- عدم توانایی اتصال مستقیم ورودی و خروجی به پردازنده ی گرافیکی

بسیاری از برنامه های واقعی نیاز دارند که حجم بسیار زیادی از داده ها را به عنوان ورودی دریافت کرده و همین طور حجم وسیعی از داده ها را به عنوان خروجی بازگردانند. در حال حاضر اطلاعات باید از ورودی ها و خروجی ها گرفته شده و در اختیار CPU قرار گیرند و سپس CPU آن ها را در اختیار پردازنده های گرافیکی قرار دهد. این انتقال اطلاعات از طریق CPU باعث افت کارایی و سرعت خواهد شد و مسلماً یکی از نقطه ضعف های سامانه های

پردازش موازی است که با کمک پردازنده های گرافیکی کار می کنند. در صورتی که بتوان معماری را به گونه ای تغییر داد که پردازنده های گرافیکی بتواند به صورت مستقیم¹ به ورودی ها و خروجی ها و یا به دیسک سخت دسترسی پیدا کند، سرعت و کارایی پردازنده های گرافیکی بهتر خواهد شد.

3- عدم پشتیبانی از حافظه های سلسله مراتبی

با توجه به اینکه کاربران مختلف نیاز به سیستم هایی با ابعاد گوناگون دارند، معمولاً سازندگان تراشه های حافظه، آن ها را طوری طراحی می کنند که امکان قرارگیری دو، چهار و یا حتی هشت تا از آن پردازنده ها بر روی یک برد ممکن باشد. این استقرار تراشه ها بدین سبک به آسانی میسر نیست. در درجه ی نخست باید یک سخت افزار مناسب برای این کار در نظر گرفته شود که بتواند از چنین سلسله مراتبی پشتیبانی کند. در درجه ی دوم باید این ترتیب قرارگیری حافظه ها از لحاظ نرم افزاری هم مورد بررسی قرار گیرد و هماهنگی های لازم برای امکان کار با چنین تراشه هایی فراهم شود. در صورتی که پردازنده های گرافیکی بتوانند از این سلسله مراتب حافظه پشتیبانی کنند، مشکل نخست نیز حل خواهد شد زیرا توسعه دهندگان سامانه ها قادر خواهند بود با توجه به نیاز سیستم برای پردازنده ی گرافیکی آن، حافظه در نظر بگیرند.

مقایسه ی نسل های مختلف پردازنده های گرافیکی

برای اینکه جمع بندی آنچه تاکنون بیان شد آسان تر گردد در جدول 0-1 جزئیات سه نسل اصلی پردازنده های گرافیکی شرکت NVIDIA را با هم مقایسه می کنیم.

جدول 0-1: پردازنده های گرافیکی نسل های مختلف شرکت NVIDIA

Fermi	GT200	G80	GPU
3000	1400	681	تعداد ترانزیستور (میلیون)
512	240	128	هسته های کودا
FMA ops/clock 256	30 FMA ops/clock	ندارد	توانایی انجام محاسبات اعشاری با دقت مضاعف

¹ از طریق DMA

توانایی محاسبات اعشاری با دقت ساده	MADS ops/clock 128	MADS ops/clock 240	MADS ops/clock 512
تعداد SFU	2	2	4
تعداد واحدهای مدیریت وارپ در هر SM	1	1	2
حافظه مشترک	16 KB	16 KB	قابل انطباق بین 16 KB یا 48 KB
کش L1	ندارد	ندارد	قابل انطباق بین 16 KB یا 48 KB
کش L2	ندارد	ندارد	768 KB
واحد ECC	ندارد	ندارد	دارد
تعداد کرنل‌ها هم‌زمان	1	1	16
پهنای باند حافظه	32 بیت	32 بیت	64 بیت

NVIDIA TESLA C2050 2.6.7

پردازنده‌های تسلا سری 20 که بر اساس معماری فرمی طراحی شده‌اند دارای بسیاری از ویژگی‌های ضروری که در سامانه‌های پردازش موازی مورد نیاز است هستند. این ویژگی‌ها شامل پشتیبانی از C++، استفاده از فناوری ECC و افزایش سرعت 7 برابری در واحد محاسبات اعشاری با دقت مضاعف نسبت به GPUهای تسلا سری 10 هستند.

جدول 2-0: ویژگی‌های پردازنده‌ی تسلا C2050 و C2070

ویژگی	توضیح
448 هسته‌ی کودا	توان ارائه‌ی 515 GFLOPS پردازش به ازای هر GPU

ECC	کنترل خطا بر روی حافظه های L1، L2، RAM و ثباتها
6 GB حافظه ی DDR5	توانایی حل بسیاری از مسائل با ابعاد وسیع
موتور GIGATHREAD	افزایش سرعت انجام برنامه ها با تغییر سریع محتوا
انتقال ناهمگام اطلاعات ¹	این GPU ها قادر هستند تا زمانی که پردازش قبلی به پایان نرسیده است اطلاعات مورد نیاز برای قسمت های بعدی را به حافظه بیاورند. به عبارت دیگر GPU می تواند در حین اجرای یک قسمت عملیات انتقال اطلاعات را نیز انجام دهد.
پشتیبانی وسیع نرم افزاری	می توان با کمک C، C++، OpenCL، Fortran و Direct Compute برای کار با تسلا سری 20 برنامه نوشت. از طرفی NVIDIA Parallel Nsight tool نیز برای کاربران Visual Studio فراهم شده است.

¹Asynchronous Data Transfer

3.

فصل سوم

معماری پاسکال

معماری پاسکال

در سال 2016 در کنفرانس خبری واقع در شهر san Jose مدل جدیدی از معماری نسل جدید پردازنده های شرکت NVidia را معرفی شد. این معماری که تحت عنوان Pascal معرفی شد با پرچم داری نظیر NVidia Tesla P100 پا به عرصه گشود. بر اساس این کنفرانس و دیتاشیتی که از این معماری برای پردازنده p100 عرضه شد این پردازنده گرافیکی بالاترین کارایی را نسبت به تمام پردازنده های گرافیکی موجود در بازار به ارمغان آورده است. همچنین استفاده آن نوع پردازنده گرافیکی توانسته است بهره ای پردازش های HPC، محاسباتی تکنیکی، یادگیری ژرف و محاسبات با حجم های بالا بر روی بستر های مختلف دیتا سنتر را به حداکثر میزان خود برساند.

در این گزارش سعی بر آن شده است که نیم نگاهی به ویژگی های این معماری تازه وارد و پردازنده گرافیکی مخصوص به آن را شرح داده شود.



3.1 Tesla P100: کارایی بی نظیر در کنار ویژگی های بارز و چشم نواز

برای پردازش های گرافیکی

پردازنده گرافیکی GP100 که از معماری Tesla P100 بهره جویی کرده است ویژگی های مخصوص به خود را دارد که این پردازنده را نسبت به سایر رقبای خود در بازار تمیز داده است. این ویژگی ها عبارت اند از:

عملکردی خارق العاده در زمینه های گوناگونی همچون HPC، یادگیری ژرف و استفاده از چندین پردازنده گرافیکی به منظور محاسبات سریع.

NVLink – مدل جدید شرکت Nvidia به منظور ارتباط سریع با پهنای باند بالا بین چنین کارت پردازنده گرافیکی. این تکنولوژی قادر است بیشترین میزان پردازنده گرافیکی را به هم متصل نماید.

HBM2 – سریعتتر، با ظرفیتی بالا همراه با قابلیت پشته برای معماری پردازنده های گرافیکی

بهینه سازی معماری یکپارچه و وقفه های بین پردازنده اصلی و گرافیکی که به طور چشمگیری به افزایش سرعت در مدل های برنامه نویسی می انجامد.

تکنولوژی 16 نانومتری FinFet که ویژگی های زیادی از جمله سرعت بالاتر، اشکال زدایی بالا، کاهش انرژی را برای پردازنده گرافیکی به ارمغان می آورد.



شکل 15: پردازنده ی تسلا P100

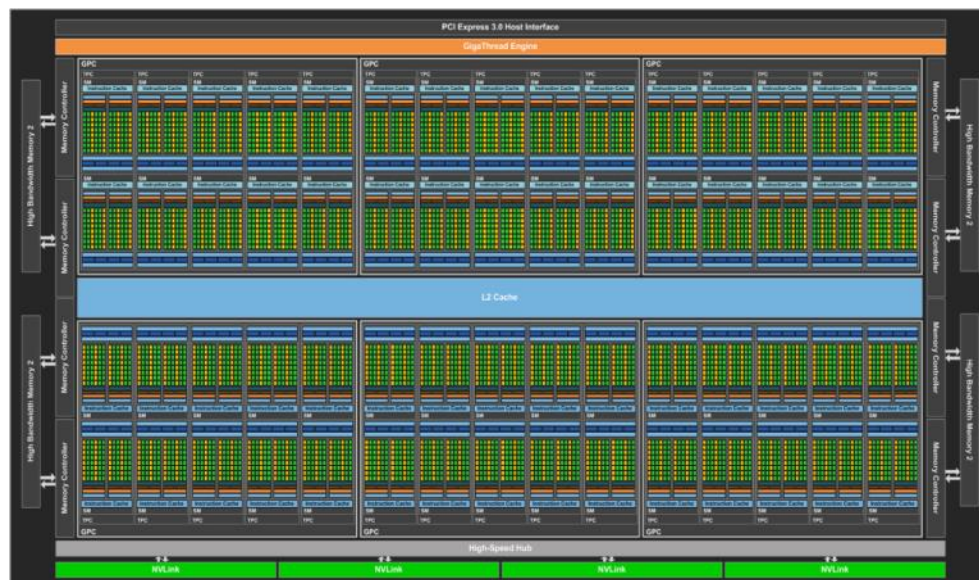
3.2 پیشنهادی عالی برای تکنولوژی NVLink و سرورهای PCI-Express

پردازنده Tesla P100 به دو صورت میتوانند با هم ارتباط برقرار کنند. یک راه سنتی که پردازنده های گرافیکی با استفاده از پورت های PCI-Express به هم و به پردازنده های گرافیکی خارج از سرور متصل خواهد شد. راه دیگری از استفاده از ماژول SXM2 برای سرور های NVLink می باشد. پردازنده ی P100 برای سرور های مبتنی بر بستر PCI-express این امکان را برای دیتاسنتر ها فراهم میکند که تمامی پردازنده های اصلی و گرافیکی با استفاده از این درگاه با هم در ارتباط باشند. همچنین P100 برای سرور های مبتنی بر مدل NVLink این امکان را فراهم میکند که برنامه هایی را به صورت توزیع شده بر روی چند پردازنده گرافیکی به منظور محاسباتی سنگین همچون یادگیری سریع فراهم آورد

3.3 معماری پاسکال GP100: از هر جهتی سریعتر

بنا به رسم قدیم شرکت Nvidia، با هر مدل جدیدی از معماری که معرفی میکند میزان مصرف انرژی را کاهش و عملکرد را افزایش میدهد. این رسم برای معماری پاسکال نیز برقرار بوده. قلب تپنده ای محاسبات پردازنده های گرافیکی را SM ها یا همان Streaming multiprocessor ها تشکیل میدهند. SM ها وظیفه ی ساخت، مدیریت، زمانبندی و اجرای عملیات های هر ترد را برعهده دارد.

همانند مدل های قدیمی پردازنده های Tesla، GP100 نیز از آرایه های کلاستر های پردازش گرافیکی (GPC)، SM ها و کنترلر های حافظه تشکیل شده است. GP100 توان عملیاتی خود را وامدار شش GPC که هر کدام از این کلاستر ها شامل 60 SM و همچنین 8 کنترلر 512 بیتی تشکیل شده است. معماری پاسکال عملکرد بهتر خود را تنها وامدار افزایش تعداد SM ها نیست بلکه در این معماری هر کدام از SM ها به صورت کارآمد تری نسبت به مدل های پیشین خود کار میکنند. برای روشن تر شدن این موضوع میتوان به این مورد اشاره کرد که هر SM شامل 64 هسته CUDA و 4 واحد بافتی برای عملیات های کار با تصویر تشکیل شده است. روی هم رفته میتوان به این نتیجه رسید که در این معماری عملا 3840 واحد CUDA و 240 واحد بافتی وجود دارد.



شکل 16: بلوک دیاگرامی از P100

بهینه سازی عملکرد پردازند گرافیکی در محاسبات و همچنین کاهش میزان انرژی به نسبت با این عملکرد دو هدف اصلی در این معماری پردازنده گرافیکی است. میزان تغییراتی که در SMهای این معماری نسبت به مدل قبلی خود یعنی معماری ماکسول صورت گرفت شایان این امر است. شاید بتوان گفت که این عملکرد مدیون نقش تکنولوژی شانزده نانومتری Fin_FET است؛ به طوری که با ارائه این تکنولوژی، به طور نسبی تمامی قطعاتی که از این تکنولوژی بهره بردند به میزان قابل قبولی عملکرد آنها افزایش یافته است.

جدول زیر یک مقایسه پیشرفته P100 با مدل های پیشین Tesla را آورده است :

Tesla Products	Tesla K40	Tesla M40	Tesla P100 (NVLink)	Tesla P100 (PCIe)
GPU / Form Factor	Kepler GK110 / PCIe	Maxwell GM200 / PCIe	Pascal GP100 / SXM2	Pascal GP100 / PCIe
SMs	15	24	56	56
TPCs	15	24	28	28
FP32 CUDA Cores / SM	192	128	64	64
FP32 CUDA Cores / GPU	2880	3072	3584	3584
FP64 CUDA Cores / SM	64	4	32	32
FP64 CUDA Cores / GPU	960	96	1792	1792
Base Clock	745 MHz	948 MHz	1328 MHz	1126 MHz
GPU Boost Clock	810/875 MHz	1114 MHz	1480 MHz	1303 MHz
FP32 GFLOPs	5040	6844	10608	9340
FP64 GFLOPs	1680	213	5304	4670

Texture Units	240	192	224	224
Memory Interface	384-bit GDDR5	384-bit GDDR5	4096-bit HBM2	3072-bit HBM2 (12GB) 4096-bit HBM2 (16GB)
Memory Bandwidth	288 GB/s	288 GB/s	732 GB/s	549 GB/s (12GB) 732 GB/s (16GB)
Memory Size	Up to 12 GB	Up to 24 GB	16 GB	12 GB or 16 GB
L2 Cache Size	1536 KB	3072 KB	4096 KB	4096 KB
Register File Size / SM	256 KB	256 KB	256 KB	256 KB
Register File Size / GPU	3840 KB	6144 KB	14336 KB	14336 KB
TDP	235 Watts	250 Watts	300 Watts	250 Watts
Transistors	7.1 billion	8 billion	15.3 billion	15.3 billion
GPU Die Size	551 mm ²	601 mm ²	610 mm ²	610 mm ²
Manufacturing Process	28-nm	28-nm	16-nm	16-nm

Pascal SM 3.4

هر واحد SM از پردازنده GP100 شامل 64 واحد single precision (FP32) است. در مقابل، در معماری ماکسول و یا کپلر که مدلی قدیمی تر نسبت به پاسکال محسوب میشود تعداد این واحدهای تک اعشاری به ترتیب 128 و 192 عدد بود. هر واحد SM پردازنده GP100 خود به دو بلاک اصلی پردازشی تقسیم میشود که هر کدام شامل 32 واحد single precision هسته کودا، بافر دستور العمل ها، سیستم مدیریتی وارپ ها و دو واحد توزیع می باشند. این درحالی است که پردازنده GP100 نصف تعداد پردازنده های کودا را نسبت به معماری ماکسول دارد. ولی با این حال تعداد و حجم رجیسترها در این معماری با معماری پیشین خود برابر است.



شکل 17: SMهای پردازنده ی تسلا P100

تعداد رجیسترهای GP100 که در یک SM با مدل های Maxwell GM200 و کپلر GK110 برابر است، ولی مدل GP100 تعداد بیشتری SM نسبت به مدل های پیشگفته دارد فلذا تعداد رجیسترهای این مدل بیشتر از سایر مدل ها خواهد بود. این امر بدین معنا است که هر ترد در پردازنده گرافیکی به تعداد بیشتری رجیستر دسترسی دارد. در یک کلام میتوان گفت که مدل GP100 تعداد بیشتری ترد، وارپ و بلاک را در مقایسه با نسل های قبل خود پشتیبانی میکند.

مقدار حافظه مشترک نیز در مدل GP100 با توجه به تعداد زیاد SM ها افزایش یافته است و نکته دیگر این که پهنای باند حافظه مشترک به طور نسبی دوبرابر نسل پیشین خود شده است. نسبت بالاتر حافظه مشترک، رجیسترها و وارپ ها در هر SM در مدل GP100 این امکان را به کاربر داده است که کد ها را به صورت بهینه به اجرا در آورد. به بیانی دیگر تعداد وارپ بیشتری برای هر کنترل کننده عملگری وجود دارد، تعدادی بیشتری عملیات load میتواند صورت بگیرد و همچنین پهنای باند برای ارتباط هر ترد با حافظه مشترک افزایش یافته است. همه اینها باعث افزایش عملکرد پردازنده گرافیکی GP100 شده است .

در مقایسه با معماری kepler، SMهای معماری پاسکال مدل ساده تری از معماری را پیش روی میگیرد به طوری که میزان فضایی که بر روی دای استفاده میکنند کمتر و به مراتب میزان میزان مصرف انرژی نیز برای انتقال داده ها در این مدل کمتر از معماری های پیش از خود خواهد بود. همچنین پاسکال یک مدل مدیریتی برای

افزایش سرعت عملیات های Load و store فراهم کرده که به مراتب عملکرد این دستور العمل ها را برای واحد های ممیز شناور افزایش داده است. سیستم مدیریت SM ها در معماری پاسکال یک سیستم هوشمند بوده به طوری که درمقایسه با معماری مکسول عملکردی بهتر را نشان میدهد.

3.5 عملکرد بالا در دقت مضاعف

به دلیل اهمیت عملیات های دقت مضاعف در محاسبات تکنیکی و کدهای HPC، یکی از هدف های اصلی برای طراحی چنین معماری ای یعنی معماری پاسکال افزایش عملکرد عملیات های دقت مضاعف بود. هر SM در پردازنده های GP100 دارای 32 واحد FP64 هستند که نسبت 2 به 1 با واحد های Single Precision دارند. این درحالی است که این نسبت در مقایسه با پردازنده Kepler GK110 3 به 1 است. فلذا این امر، این امکان را به پردازنده گرافیکی P100 میدهد که عملیات های دقت مضاعف را با استفاده از واحدهای FP64 در حجم کار بالا بهتر کند.

همانند مدل های پیشین پردازنده های گرافیکی، پردازنده GP100 نیز قوانین IEEE 754-2008 را چه برای مدل دقت تکی چه دقت مضاعف پشتیبانی میکند.

پشتیبانی از عملیات ریاضی FP16 برای تسریع در برنامه های یادگیری ژرف:

رشد روز افزون رشته هایی از قبیل یادگیری ژرف این خواسته را به وجود آورده است که یک وسیله محاسباتی سریع به خصوص برای این زمینه فراهم آید. یادگیری ژرف ثابت کرده است که میتواند با دقتی بالا و وفق پذیری مناسبی برنامه های کاربری از قبیل توصیف خودکار تصاویر، خودرو های بدون سرنشین، یادگیری زبان های طبیعی و ترجمه ی آنها و از همه مهم تر تولید خودکار هنر های کامپیوتری میتوان نام برد.

برخلاف دیگر برنامه های محاسباتی تکنیکی که نیازمند عملکرد سریع واحدهای اعشاری است، معماری شبکه های عمیق به دلیل متدهای backpropagation خطاهایی به طور طبیعی در حین مسیر یادگیری خود دارند. ذخیره سازی داده های FP16 در مقایسه با داده های دقت بالاتری نظیر FP32 و FP64 به نسبت فضای کمتر از حافظه شبکه عصبی مصرف میکند، که این امر این امکان را به ما میدهد تا شبکه عصبی بزرگ تری هنگام آموزش شبکه و یا اجرای شبکه طراحی کنیم. استفاده از FP16 عملکرد شبکه عصبی را تا دو برابر نسبت به استفاده از عملیات FP32 بالا برده است و به طور مشابه استفاده از FP16 در انتقال داده ها زمان کمتری نسبت به عملیات های FP32 و FP64 خواهد داشت.

سیستم GP100 SM ISA این امکان را فراهم کرده است که عملیات محاسباتی FP16 را به صورت دوتایی در یک زمان در یک هسته کودا اجرا نماید. گفتنی است که ریجسترهای 32 بیتی GP100 میتواند دو مقدار FP16 را در خود جای دهد.

3.6 افزایش سرعت عملکرد عملگرهای Atomic

عملیاتهای اتمیک یکی از مهمترین عملیات در برنامه نویسی به صورت موازی است. چرا که به برنامه نویس این امکان را می دهد که چند ترد که به صورت موازی به اجرا در آمده اند عملیاتهای خواندن-تغییر-نوشتن را در یک حافظه مشخص به درستی انجام دهند. معماری کپلر به طور چشمگیری این قابلیت را بهبود بخشید به طوری که این امکان نه تنها بر روی حافظه های محلی میسر بود بلکه میتوان بر روی حافظه های اصلی یا همان global نیز عملیات اتمیک را به اجرا در آورد. اگر چه هر دو معماری فرمی و کپلر هزینه زیادی برای عملیاتهای اتمیک قائل بودند.

با روی کار آمدن معماری مکسول، بهبود مناسبی برای متغیرهای 32 بیتی صحیح فراهم شد به طوری که این عملیات اتمیک به صورت سخت افزاری و کاملاً محلی پشتیبانی میشد. همچنین این امکان سخت افزاری برای حالت های 32 بیتی و 64 بیتی برای عملیات مقایسه-تعویض (CAS) برقرار است که این سربار را از برنامه کم خواهد کرد. (درمقایسه با معماری های فرمی و کپلر که این عملیات به صورت نرم افزاری کنترل میشدند).

حال پردازنده GP100 با بهبود بخشی به عملیاتهای اتمیک این امکان جمع FP64 را در حافظه global به صورت سخت افزاری فراهم کرده است. تابع atomicAdd در پلتفرم کودا هم در اعداد صحیح و اعشاری 32 بیتی هم 64 بیتی را فراهم کرده است. در مدل های قبل برای جمع های FP64 از حلقه های CAS استفاده میشد که به طور کلی بسیار کند تر مدل های سخت افزاری و محلی است.

قابلیت محاسباتی 6.0

GP100 قابلیت محاسباتی جدید 6.0 را پشتیبانی میکند. جدول زیر تفاوتهای پارامترهای مختلف قابلیت های محاسباتی را نشان می دهد.

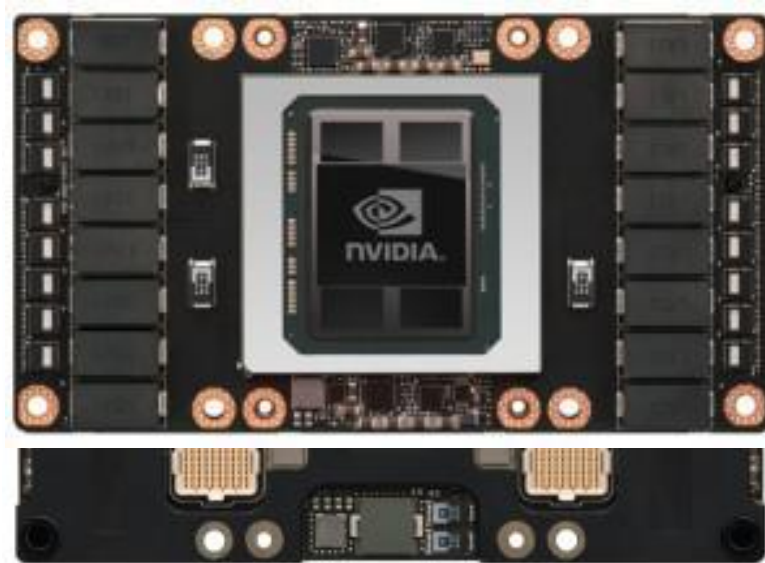
GPU	Kepler GK110	Maxwell GM200	Pascal GP100
Compute Capability	3.5	5.2	6.0

Threads / Warp	32	32	32
Max Warps / 64 Multiprocessor		64	64
Max Threads / 2048 Multiprocessor		2048	2048
Max Thread Blocks / 16 Multiprocessor		32	32
Max 32-bit Registers / 65536 SM		65536	65536
Max Registers / Block	65536	32768	65536
Max Registers / Thread	255	255	255
Max Thread Block Size	1024	1024	1024
CUDA Cores / SM	192	128	64
Shared Memory Size / 16K/32K/48K SM Configurations (bytes)		96K	

3.7 بهبودی شگرف در حافظه

بسیاری از برنامه های امروزی به دلیل محدودیت های پهنای باند حافظه دچار مشکل هستند مخصوصاً برنامه ها مرتبط با HPC. امروزه، برنامه نویسان برنامه های پردازش سریع بدون توجه به پردازنده ای که با آن کار میکنند، تلاش بسیاری برای بهبود عملکرد برنامه بلاخص برای دسترسی به حافظه دارند. این تلاش ها ممکن است در حیطه نزدیک نگاه داشتن واحد های پردازشی به سلسله مراتب های حافظه خواهد بود. برخی از برنامه ها از قبیل برنامه های یادگیری سریع که خود شامل چنین لایه از شبکه عصبی است، بایستی با استفاده از یک دیتابیس وسیع و حجیم آموزش داده شود، بزرگترین محدودیتی که دارد، محدودیت حافظه است. بنابراین حافظه ها دو محدودیت بزرگ برای عملیات های محاسباتی دارند: پهنای باند و ظرفیت

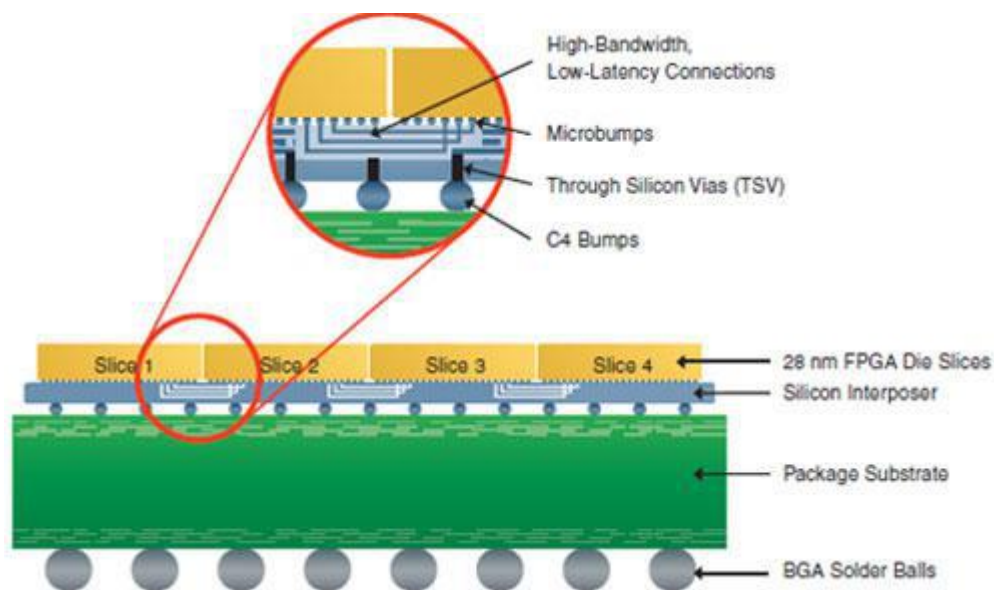
تسلا P100 هر دوی این محدودیت ها را با استفاده از تکنولوژی حافظه های پشته ای برطرف کرده است. تکنولوژی ای که چندین لایه از واحد های DRAM را قادر میسازد تا به صورت عمودی در یک پکیج درون پردازنده گرافیکی با هم در ارتباط باشند. تسلا P100 اولین پردازنده گرافیکی است که با استفاده از تکنولوژی پهنای باند سریع 2 (HBM2) عرضه شده است. این تکنولوژی پهنای باند وسیع و همچنین ظرفیتی بالا، چیزی در حدود 2 برابر، و همچنین بهبود مصرف انرژی را نسبت به GDDR5 دارد.



شکل 18: شتابدهنده ی P100 از پشت

بر خلاف GDDR5 که شامل تکه های جدا از چپ های حافظه پیرامون پردازنده گرافیکی بود، در تکنولوژی HBM2 شامل یک یا چند پشته عمودی در چند دای مموری خواهد بود. دای های مموری با استفاده از سیم های ریزی که درون-سیکِلونی¹ و microbump ها مرتبط هستند.

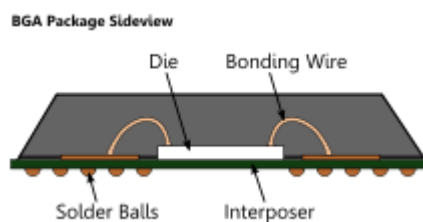
¹Through-silicon



شکل 19: نمایشی از شمای دای و میکرو بامپ ها

هر دای 8 گیگابایت HBM2، شامل 5000 سوراخ های درون-سیلیکان می باشد. در نهایت یک اینترپوزر¹ passive از نوع سیلیکان برای مرتبط کردن پشته های حافظه و واحد پردازنده گرافیکی مورد استفاده قرار میگیرد. ترکیب این HBM2 ها، دای پردازنده گرافیکی و اینترپوزهای سیلیکانی در یک پکیج 55 در 55 میلیمتری تحت عنوان BGA package عرضه می شود.

تسلا P100 شامل 4 دای از پشته های HBM2 و 16 گیگابایت از حافظه است. این پردازنده قادر است با پهنای باند 720 GB/s ارتباط برقرار کند (گفتنی است که این پهنای باند در حالت پیک می باشد) به طوری که این پهنای باند سه برابر بیشتر از مدل تسلا M40 است.



شکل 20: محل اینترپوزر

¹ Interposer

ECC Memory 3.8

از دیگر مزیت های تکنولوژی HBM2 میتوان به پشتیبانی محلی از رفع خطاهای کد^۱ اشاره کرد. این ویژگی قابل اطمینان بودن عملیات های محاسباتی برای برنامه هایی که به خرابی داده ها حساس هستند فراهم میکند. به عنوان مثال میتوان به برنامه هایی که در کلاستر های عظیم و یا سوپر کامپیوتر ها که پردازنده گرافیکی حجم بالایی از از دیتاست ها را پردازش میکنند اشاره کرد.

تکنولوژی ECC یک بیت از خطا را قبل از آن که به سیستم آسیبی بزند شناسایی و رفع می کند. در مدل GDDR5 تکنولوژی ECC را برای محافظت از محتوای حافظه به صورت محلی پشتیبانی نمیکند، این رفع یابی بایستی در درگاه داده صورت پذیرد که خود موجب به کاهش سرعت می شود.

در مدل پردازنده گرافیکی Kepler GK110 امکان حفاظت از داده های حافظه با استفاده از تکنولوژی ECC به صورت محدود می دهد. این محدودیت بدین معناست که جزئی از حافظه از این قابلیت استفاده میکند. 6.25% از کل حافظه GDDR5 به منظور بررسی ECC رزور شده است. در حالتی که برای پردازنده K40 حافظه 12 گیگابایتی را در نظر بگیریم، 750 مگابایت تنها سهم یان تکنولوژی شده و مابقی حافظه یعنی 11.25 گیگابایت بایستی بر روی درگاه داده عیب یابی شوند. همچنین دسترسی به بیت های ECC باعث کاهش خفیفی در پهنای باند حافظه نسبت به عدم استفاده از ECC خواهد بود. بنابراین HBM2 به دلیل آن که ECC را به صورت محلی و سخت افزاری پشتیبانی میکند، دیگر این کاهش پهنای باند برای این معماری بی تاثیر خواهد بود. همچنین ECC در هر زمانی فعال خواهد بود. همانند پردازنده گرافیکی GK110، پردازنده GP100 نیز از رفع خطای تکی^۲ و شناسای خطای دوتایی^۳ (SECCDED) پشتیبانی میکند.

NVLink 3.9 برای ارتباط سریع

NVLink یک تکنولوژی جدید و سریع به منظور ارتباطات بین پردازنده های گرافیکی است که توسط شرکت Nvidia معرفی شده است. با توجه به پشتیبانی ای که پردازنده P100 از SXM-2 میکند، تکنولوژی NVLink به این امکان را دارد که پردازنده های گرافیکی با هم و یا با پردازنده اصلی به صورت بسیار کارآمد ارتباط برقرار کنند.

¹Error correct code (ECC)

²Single error correct

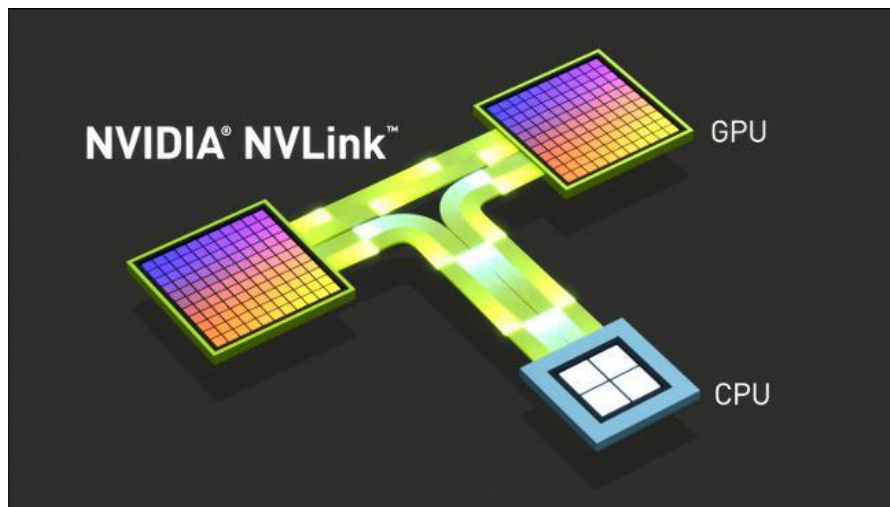
³Double error correct

امروزه، استفاده از چند پردازنده گرافیکی در یک ایستگاه‌های پردازشی^۱ همانند استفاده از چندین node در کلاستر های HPC و یا سیستم های آموزش به یادگیری های ژرف مر سوم است. ارتباطات قوی در یک سیستم چند پردازنده ای بسیار ضروری است. هدف از معرفی تکنولوژی NVLink ارائه یک سیستم ارتباطی بین پردازنده های گرافیکی است که بتواند با پهنای باند بالا با استفاده از درگاه های PCI-Express نسل سوم با هم در ارتباط باشند. از دیگر اهداف این تکنولوژی سازگاری آن با GPU ISA به منظور پشتیبانی از حافظه های محلی پردازنده های گرافیکی در سیستم های چند پردازنده ای است.

پشتیبانی از GPU ISA بدین معنا است که بتوان با استفاده از تکنولوژی NVLink پردازنده های گرافیکی را به هم متصل کند و برنامه ما به طور مستقیم بر روی داده هایی که در چند پردازنده گرافیکی قرار دارد اجرا شود. پردازنده گرافیکی همچنین میتواند عملیات اتومیک حافظه ای را بر روی پردازنده های گرافیکی کنترلی از راه دور را بر عهده گیرد که این امر باعث بهبود در توزیع مناسب برنامه ها می شود.

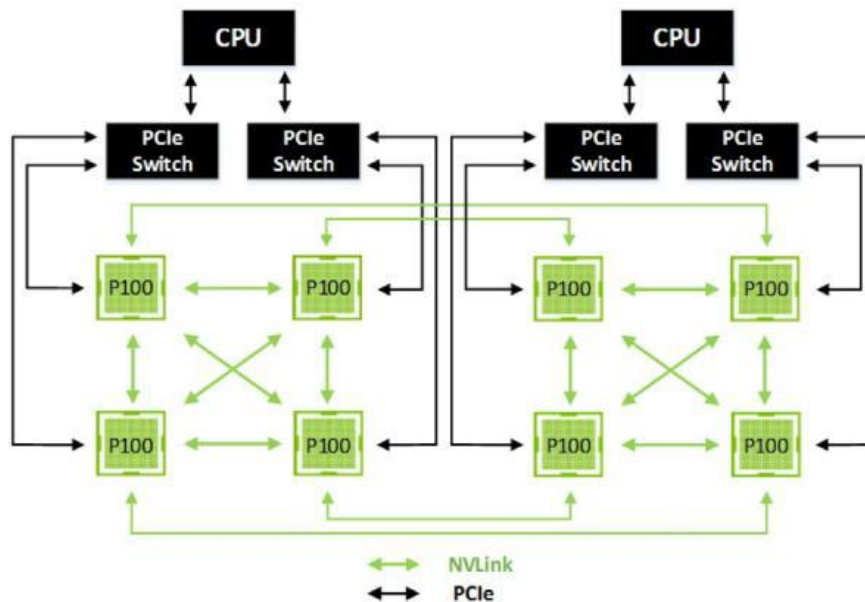
NVLink از ارتباطات سریع سیگنالی انویدیا^۲ بهره می جوید. NVHS داده ها را بر روی جفت های تفاضلی که با سرعتی حداکثر 20 گیگابیت بر ثانیه دارد انتقال میدهد. هشت عدد از این جفت های تفاضلی یک sub-link را تشکیل میدهد که داده را در به صورت یک طرفه انتقال میدهد و دو sub-link - هر کدام برای یک جهت - یک لینک را تشکیل میدهد که این لینک دو پردازنده گرافیکی یا یک پردازنده گرافیکی با یک پردازنده اصلی وصل میکند. یک لینک می تواند تا 40 گیگابیت بر ثانیه را به صورت دو طرفه انتقال دهد. چند لینک میتواند با هم متصل شوند و گنگ^۳ را تشکیل دهد که پهنای باند وسیع تری را برای ارتباط بین پردازنده ها فراهم آورد. NVLink ای که در پردازنده P100 به کار رفته است از چهار لینک تشکیل شده است. این بدان معنا است که بیشینه پهنای باندی که به صورت تئوری قابل انتقال به صورت دوطرفه خواهد بود برابر است با 160 گیگابیت در ثانیه.

¹ Workstation² Nvidia's High-speed Signaling interconnect (NVHS)³ gangs



شکل 21: NVLink

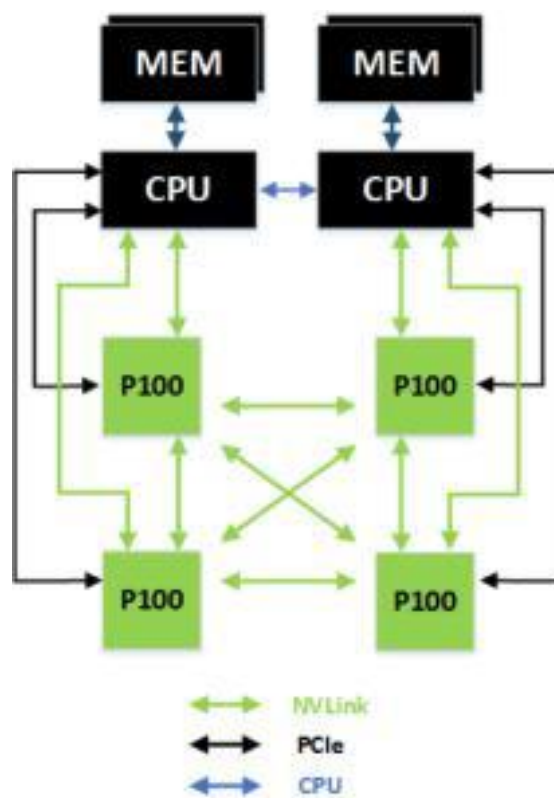
تصویر زیر 8 پردازنده گرافیکی به صورت مش های مکعبی ترکیبی^۱ که 4 تا از آنها به صورت کامل با استفاده از NVLink به هم مرتبط هستند و یک NVLink نیز بین این دسته های چهارتایی متصل است نشان میدهد. هر پردازنده گرافیکی در دسته چهارتایی خود به پردازنده اصلی دسته خود به طور مستقیم با استفاده از درگاه PCI-E متصل است.



شکل 22: هشت پردازنده گرافیکی که به صورت hyper cube قرار دارند

^۱ Hybrid cube mesh

در حالی که وظیفه اصلی NVLink وصل کردن چندین پردازنده گرافیکی GP100 به یکدیگر است، این امکان را که پردازنده گرافیکی GP100 را به پردازنده های اصلی IBM متصل نمود، وجود دارد. تصویر زیر معرف این موضوع است که چهار پردازنده گرافیکی با استفاده از تکنولوژی NVLink به CPUها متصل شده است. در این پیکربندی، هر پردازنده گرافیکی مجموعاً دارای سرعت 120 گیگابایت بر ثانیه به صورت دو طرفه با 3 پردازنده گرافیکی دیگر است و 40 گیگابایت ارتباط دوطرفه با پردازنده اصلی است.



شکل 23: 4 پردازنده گرافیکی که با استفاده از تکنولوژی NVLink به پردازنده اصلی و خود متصل اند.

3.10 افزایش بهره وری برنامه نویسی با استفاده از حافظه یکپارچه¹

حافظه یکپارچه یکی از اصلی ترین ویژگی های برنامه نویسی به زبان کودا است به طوریکه عملیات برنامه نویسی را به طور چشمگیری برا برنامه نویس راحت میکند. این بدان معناست که میتوان برنامه و اطلاعات لازم را به

¹ Unified memory

پردازنده گرافیکی فرستاد ولی در عین حال از یک حافظه به منظور دسترسی پردازنده گرافیکی و پردازنده اصلی به حافظه استفاده نمود. مدل پاسکال GP100 با ارتقای سطح عملکرد حافظه یکپارچه پیشرفت شگرفی در این امر انجام داده است.

با معرفی کودا 6، حافظه های یکپارچه نیز معرفی شده که در آن یک حافظه که هم از طریق پردازنده گرافیکی و هم پردازنده اصلی قابل دسترسی است. این حافظه ها از دو پردازنده اصلی و گرافیکی با استفاده از یک اشاره گر قابل دسترسی بود. سیستم کودا به صورت خودکار داده هایی که در حافظه یکپارچه قرار دارد را بنا به نیاز به حافظه اصلی یا گرافیکی منتقل میکند، بدین صورت که گویی پردازنده اصلی برنامه ای را در حافظه اصلی خودش پردازش میکند و یا پردازشگر گرافیکی داده را در حافظه گرافیکی خود دارد.

پلتفرم کودا 6 برای حافظه های یکپارچه با محدودیت هایی روبه رو بود. یکی از این محدودیت ها آن بود که اگر سیستم کنترل حافظه تماسی با پردازنده اصلی داشته باشد باید عملیات همزمان سازی قبل از اجرای کرنل ها صورت پذیرد. همچنین امکان دسترسی همزمان پردازنده گرافیکی و پردازنده اصلی به یک نقطه از حافظه امکان پذیر نبود. از دیگر محدودیت ها نیز میتوان به محدود بودن حافظه یکپارچه به اندازه حافظه پردازنده گرافیکی اشاره کرد.

CPU Code

```
void sortfile(FILE *fp, int N) {
    char *data;
    data = (char *)malloc(N);

    fread(data, 1, N, fp);

    qsort(data, N, 1, compare);

    use_data(data);

    free(data);
}
```

CUDA 6 Code with Unified Memory

```
void sortfile(FILE *fp, int N) {
    char *data;
    cudaMallocManaged(&data, N);

    fread(data, 1, N, fp);

    qsort<<<...>>>(data, N, 1, compare);
    cudaDeviceSynchronize();

    use_data(data);

    cudaFree(data);
}
```

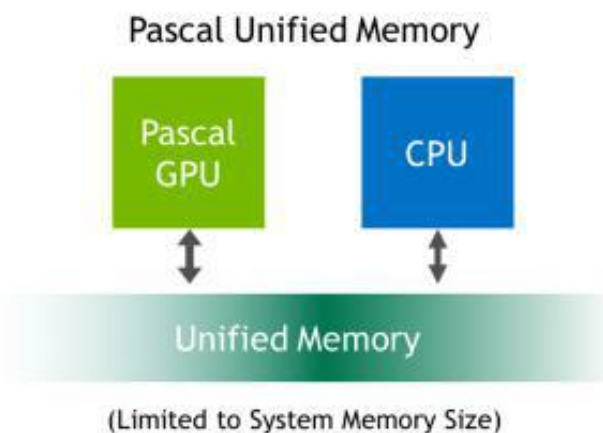
شکل 24: نمونه کد کودا 6

3.10.1 حافظه یکپارچه پردازنده P100

با توجه به مزیت هایی که حافظه یکپارچه کودا 6 برای کاربر به ارمقان آورد، در مدل پاسکال GP100 تعداد بیشتری از ویژگی برای راحتی برنامه نویسان و ارتباطات بین حافظه های پردازنده های گرافیکی و اصلی نمود پیدا

کرد، به طوری که ارتباط CPU برای اجرای برنامه های موازی شده با پردازنده گرافیکی با سرعت بالایی ارتقا پیدا کرد. دو ویژگی سخت افزاری مهم عامل این پیشرفت ها بوده اند: پشتیبانی از آدرس دهی به میزان بالایی از حافظه و رفع یابی آدرس دهی صفحه ای.

GP100 امکان آدرس دهی GPU را تا 49 بیت برای آدرس دهی مجازی فعال کرد. این میزان عظیم از آدرس دهی به منظور پوشاندن حافظه های اصلی نسل جدید، همپا با پردازنده های گرافیکی جدید لازم بود. بنابراین، حافظه یکپارچه GP100 قابلیت دسترسی به تمام آدرس های پردازنده گرافیکی و پردازنده اصلی دارد.



پشتیبانی پردازنده گرافیکی از خطایابی آدرس دهی صفحه یکی از بنیادی ترین تغییراتی است که در این نسل از پردازنده گرافیکی اتفاق افتاده است. ترکیب ویژگی های آدرس دهی مجازی 49 بیتی و عیبابی آدرس دهی صفحه مزایای فراوانی را برای کاربران فراهم کرده است. اول از همه بایستی گفت که عیبابی صفحه بدین معناست که پلتفرم کودا دیگر نیازی به همزمان سازی سیستم مدیریت حافظه ای قبل از اجرای کرنل ها ندارد. اگر کرنلی که بر روی پردازنده گرافیکی در حال اجراست به حافظه ای از پردازنده گرافیکی دسترسی پیدا کند که درون حافظه نباشد، خطا رخ خواهد داد و این امکان فراهم می شود که صفحه مورد نظر به صورت خودکار به درون حافظه پردازنده گرافیکی انتقال یابد. متناوباً، ممکن است که صفحه مورد نظر به درون حافظه پردازنده گرافیکی نگاشته شود به طوری که ارتباط پردازنده گرافیکی با داده مورد نظر از طریق درگاه PCI-Express و یا NVLink صورت پذیرد. گفتنی است که به طور کلی نگاشت¹ سرعت عمل بالاتری نسبت به انتقال داده درون پردازنده گرافیکی دارد.

فلذا با توجه به این سیستم خطایابی صفحه در نسل جدید پردازنده های گرافیکی و کودا 8، وابستگی داده های global و حافظه ی یکپارچه تضمین خواهد شد. این بدان معناست که با پردازنده GP100، میتوان به صورت همزمان پردازنده های گرافیکی و پردازنده های اصلی به حافظه یکپارچه دسترسی داشته باشند. گفتنی است که این

¹mapping

امر در مدل های کیپر و مکسول قابل قبول نبود. به این دلیل که هیچ تضمینی وجود نداشت، زمانی که پردازنده گرافیکی در حال اجرای کرنلی باشد، پردازنده اصلی قادر به دسترسی به حافظه یکپارچه داشته باشد.

در نهایت، با پشتیبانی از این ویژگی، هر نوع اختصاص دهی حافظه توسط سیستم عامل از قبیل malloc و یا new هم از طریق پردازنده گرافیکی و هم از طریق پردازنده اصلی توسط یک اشاره گر قابل دسترسی است. (می توانید به عکس زیر مراجعه کنید. در این سیستم، استفاده از حافظه یکپارچه به صورت پیش فرض مورد استفاده است. با این حساب دیگر نیازی به یک تخصیص دهنده خاص برای استفاده تخصیص حافظه نخواهد بود.

CPU Code

```
void sortfile(FILE *fp, int N) {
    char *data;
    data = (char *)malloc(N);

    fread(data, 1, N, fp);

    qsort(data, N, 1, compare);

    use_data(data);

    free(data);
}
```

Pascal Unified Memory*

```
void sortfile(FILE *fp, int N) {
    char *data;
    data = (char *)malloc(N);

    fread(data, 1, N, fp);

    qsort<<<...>>>(data, N, 1, compare);
    cudaDeviceSynchronize();

    use_data(data);

    free(data);
}
*with operating system support
```

این نکته نیز نباید فراموش شود که چندین دستورعمل خاص به منظور فعال سازی این متد بر روی سیستم عامل نیازمند است که شرکت NVidia با همکاری Red Hat و انجمن های لینوکس این عملگر قوی را فعال کرده اند.

3.11 NVidia DGX-1 یادگیری ژرف با استفاده از سوپر کامپیوترها

دانشمندی که در حوزه پردازش اطلاعات و یا هوش مصنوعی کار میکنند، به دقت بالا، سادگی و در عین حال سرعت بالا برای سیستم های یادگیری ژرف نیازمندند. آموزش هر چه سریع تر باشد، نمو خلاقیت های پیرو این حیطه از علم و همچنین فروش محصولات وابسته به آن در بازار جهانی سریع تر خواهد بود. پردازشگر NVidia DGX-1 اولین سیستمی است که به منظور پردازش های ژرف طراحی شده است، به طوری که یک پک کامل از سخت افزار ها و نرم افزار های لازم به منظور اجرای سریع و آسان الگوریتم های لازم را فراهم آورد. این

سوپر کامپیوتر انقلابی در سیستم یادگیری ژرف محسوب میشود به طوری که عملکرد سیستم تا 170 TFLOP/s برای عملیات های FP16 ارتقا پیدا کرد.

NVidia DGX-1 اولین سروری است که از پردازنده گرافیکی تسلا P100 که استفاده کرده است. این سرور قادر است که 8 پردازنده گرافیکی P100 را پیکربندی کند.



شکل 25: یک سوپر کامپیوتر DGX-1

پیکربندی 8 پردازنده گرافیکی P100 با استفاده از تکنولوژی NVLink این امکان را فراهم می سازد که این 8 پردازشگر را به صورت 4 تایی به صورت کامل به هم مرتبط کند. بدین معنا که 4 پردازشگر به صورت دوبه دو به هم متصل هستند. هر پردازشگر در این سیستم چهار تایی با استفاده از درگاه PCI-E متصل بوده و از آن به پردازشگر اصلی متصل خواهد بود. گفتنی است که در بدترین حالت، استفاده از سرور NVidia DGX-1 در حالتی که از 8 پردازشگر P100 استفاده کند، دوازده برابر افزایش سرعت در حالتی از نسخه قبلی پردازشگر ها استفاده شود، دارد.

	DUAL XEON	DGX-1
FLOPS (CPU + GPU)	3 TF	170 TF
AGGREGATE NODE BW	76 GB/s	768 GB/s
ALEXNET TRAIN TIME	150 HOURS	2 HOURS
# NODES FOR 2HR TAT	>250*	1

4.

فصل چهارم

تسريع الگوريتم هاى يادگيرى ژرف و هوش مصنوعى با
استفاده از پردازنده گرافيكى

تسریع الگوریتم های یادگیری ژرف و هوش مصنوعی با استفاده از پردازنده گرافیکی

هدف اصلی از علم هایی نظیر یادگیری ژرف و هوش مصنوعی بر آورده کردن نیاز هایی از قبیل ساخت یک ماشین که بتواند به صورت هوشمندانه به اسان خدمت کند. این هوشمندی بدین صورت است که نیازی به یک دستور عمل خاص و یا الگوریتم های مشخص نداشته و ماشین به صورت خودکار با یادگیری از پیرامون خود نیازهای بشری را ارضا نماید. یادگیری ژرف یکی از مهمترین دست آوردهای جدید هوش مصنوعی است که در سال های اخیر رونق گرفته است.

یادگیری ژرف این امکان به مغز، خواه انسان و یا خواه ماشین، می دهد که با یادگیری از محیط پیرامون، تصمیماتی درست و به جا در موقعیت های خاص بگیرد. اما اذعان به این نکته ضروری است که آموزش به چنین ماشینی نیاز به حجم عظیمی از داده دارد. به علاوه، در مدل های یادگیری ژرف که در آن ها از لایه های زیادی استفاده شده است، پردازش این اطلاعات را برای کاربر سخت و طاقت فرسا میکند. در سال 2012، شرکت گوگل، پروژه ای تحت عنوان google's Deep Learning را که تحت مغز گوگل شناخته میشود راه انداز کرد. این پروژه به منظور شناسایی گربه ها با استفاده از فیلم های یوتوب کلید زده شد. اما این مغز مصنوعی نیازمند حداقل 2000 پردازنده اصلی و یا به عبارتی دیگر نیازمند 16000 هسته CPU، به منظور قدرت دادن به این پردازش ها بود. کمتر کمپانی ای در دنیا وجود دارد که بتواند چنین حجم عظیمی از منابع و سخت افزار را یکجا برای انجام چنین پردازشی فراهم آورد. در همین هنگام، شرکت NVidia با همکاری با دانشگاه Stanford تصمیم به استفاده از پردازنده های گرافیکی برای پردازش های مورد نیاز الگوریتم های یادگیری ژرف گرفتند. نتیجه ی چنین تصمیمی، چیزی نشد جز آن که 12 پردازنده گرافیکی از شرکت NVidia توانایی یادگیری ژرف همگام با 2000 پردازنده اصلی را برآورد کرد. همانطور که از اعداد مشخص است این کاهش میزان پردازنده، نه تنها در میزان مصرف انرژی صرفه جویی بسزایی کرد، بلکه این امکان را به دیگر شرکت ها داد تا با تهیه تعدادی که پردازنده گرافیکی کمتر از تعداد انگشتان دست، به رقابت با پردازنده های Intel در الگوریتم های یادگیری ژرف بپردازند.

خیلی ها بر این باورند که انقلاب یادگیری ژرف به سال 2012 که در این سال مسابقه ImageNet با حضور سردم داران این علم یعنی Krizhevsky, sutskever و Hilton برگزار شد. این اشخاص با به کار

گیری پردازنده های گرافیکی، الگوریتم CNN را که بدان تحت عنوان AlexNet نیز یاد میشود را راه اندازی کردند و نتیجه چیزی نشد جز آن که عملکردی دوجندان بهتر نسبت به رقبایی که از کلاستر های قوی برای این مسابقه استفاده کرده بودند، به دست آورد. لازم به ذکر است که کریژوسکی و تیمش هیچ کدی در زمینه پردازش تصویر ننوشتند. بلکه متد آنها بدین صورت بود که با استفاده از یادگیری ژرف این امکان را با ماشین دادند تا خودش شروع به یادگیری کند. آنها با طراحی یک شبکه عصبی موسوم به AlexNet و آموزش آن با استفاده از میلیون ها مثال که خود نیازمند تریلیون ها محاسبات ریاضیست این مهم را به انجام رساندند. شایان ذکر است که شبکه عصبی AlexNet بهترین کد به دست انسان نوشته شده در آن سال شد.

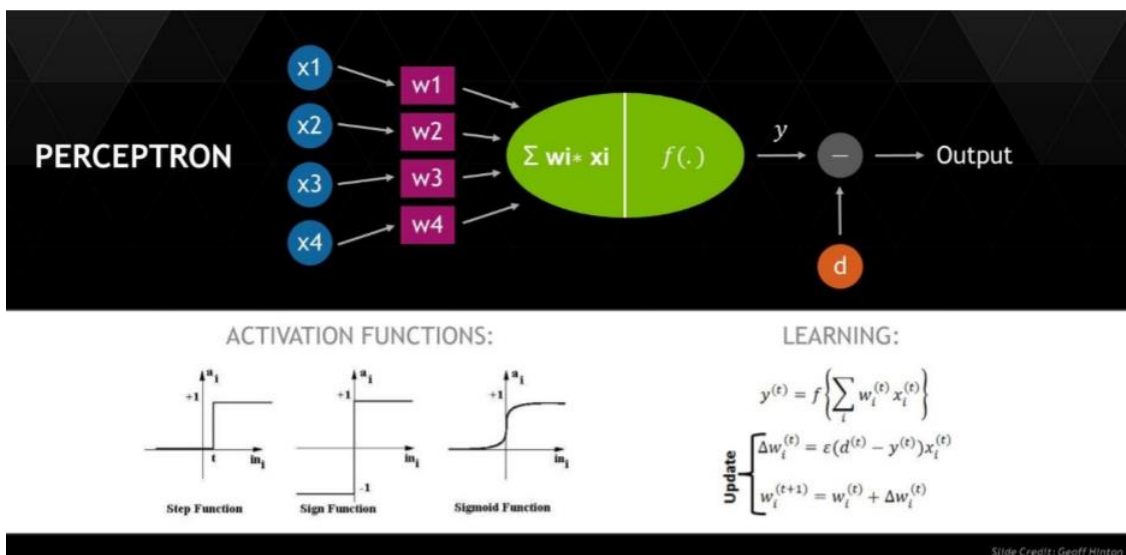
پس از آن، رفته رفته، شبکه های ژرف عصبی بر روی پردازنده های گرافیکی انقلابی را در زمینه های پردازش تصویر به طور خاص و به طور کلی الگوریتم های وابسته به یادگیری راه انداختند. کاربرد این علم به مراتب فراتر از درک ماست، به عنوان مثال استفاده از ماشین هایی که بدون سرنشین قادر به رانندگی بر روی هر سطحی و تحت هر شرایطی هستند و یا مترجم هایی که بلادرنگ هر زبانی هر با زبان دلخواه ما ترجمه سلیس میکند. امروزه، استفاده از یادگیری ماشین در کنار پردازنده های گرافیکی تمام امکانات فوق الذکر را برای ما فراهم کرده است. به طوری که استفاده از پردازنده های گرافیکی به طور چشمگیری مرحله زمان آموزش شبکه عصبی را کاهش داده است.

4.1 یادگیری ژرف در نگاه کلی

یادگیری ژرف یک تکنیکی است که مغز اسنان را با استفاده از شبکه ها عصبی و متدهای یادگیری عصبی شبیه سازی می کند. با یادگیری بیشتر، ماشین باهوش تر خواهد شد و در نتیجه جواب های دقیق تری را برای ما فراهم خواهد کرد. به عنوان مثال، کودکی را در نظر بگیرید که توسط پدر و مادرش در حال یادگیری مطالبی از قبیل شناسایی دقیق اشکال هندسی است؛ به تدریج این امکان برای کودک فراهم خواهد شد که اشکال هندسی مختلفی را بدون نیاز به والدین خود شناسایی کند. به طور مشابه، یادگیری ژرف و یا یادگیری عصبی این امکان را به سیستم میدهد که اشیا را از تصاویر شناسایی و آن ها را دسته بندی کند. این قدرت میتواند فراتر از مثال های ساده رود به طوری که سیستم هر شکلی خواه کامل باشد یا ناقص خواه واضح باشد یا پوشیده شناسایی نماید.

در سطحی ساده و مبتدی، عصب ها در مغز انسان به صورت یک هاب پذیرای ورودی های مختلف است، در نظر بگیرید. و مرتبه اهمیت هر کدام از این ورودی ها به مراتب مشخص شده است. در نهایت نیز خروجی یک عصب نیز به عنوان ورودی به عصبی دیگر منتقل خواهد شد تا تصمیمات لازم انجام شود.

عصب، همانطور که در شکل زیر مشخص است، یک مدل بسیار ساده شده از مغز انسان است. با توجه به شکل، هر پرسپترون¹ چندین ورودی که هر کدام مشخص کننده یک ویژگی از داده ورودی است، وجود دارد. وظیفه نهایی این سلسله پرسپترون ها آموختن و شناسایی و طبقه بندی سازی هر کدام از داده های ورودی ای به هنگام تست الگوریتم میباشد. گفتنی است که هر کدام از این ویژگی ها که به ورودی پرسپترون داده شده است یک وزن به خصوص که میزان اهمیت آن ویژگی را مشخص میکند نشان میدهد.



شکل 26: یک پرسپترون که مدلی ساده از شبکه عصبی است

برای مثال، یک پرسپترون را در نظر بگیرید که هدف آن شناسایی شماره صفر در عکس های دست نوشته است. واضحاً، عدد صفر به صورت های گوناگونی میتواند نوشته شود که این صور به دست خط نگارنده وابسته است. حال پرسپترون در این مرحله عکسی که شماره صفر بر آن است را به عنوان ورودی میگیرد. سپس آن را به بخش های مختلفی تقسیم میکند. به عنوان مثال این عکس میتواند به 4 بخش که مشخص کننده 4 ویژگی است تقسیم می شود. منحنی قسمت راست بالایی به ویژگی شماره 1

نگاشت می‌وشد منحنی پایین نیز به ویژگی دوم نگاشت خواهد شد و این روند به همین صورت ادامه پیدا میکند. وزنی که به هر کدام از این ویژگی‌ها تخصیص می‌یابد مشخص کننده میزان اهمیت هر کدام از این ویژگی‌هاست. قسمت بیضی شکل سبز رنگ در شکل بالا مشخص کننده مرحله است که پرسپترون هر کدام از وزن ویژگی‌ها را متناسب با ورودی آن مورد محاسبه قرار می‌دهد که در نهایت تشخیص دهد آیا عدد صفر است یا خیر.

مهم ترین ویژگی شبکه عصبی آموزش شبکه برای رسیدن به جوابی بهتر است. هر پرسپترون برای شناسایی عدد صفر از عکس، در ابتدا به صورت تصادفی مقدار دهی میشوند و سپس با متدهای یادگیری شروع به ارتقای ماتریس وزن (که ارزش های هر کدام از ویژگی های تصویر را مشخص می کند) میکند. این پرسپترون که برای عدد صفر تخصیص داده شده است، روند را ادامه می‌دهد تا با این نتیجه برسد که آیا عکس عدد صفر را نشان می‌دهد یا خیر. این فاز برنامه را نشر رو به جلو^۱ می‌گویند. چنانچه شبکه عصبی قادر به شناسایی درست عدد نباشد، باستی علت این عدم توانایی شناسایی شود، این شناسایی عدم توانایی معمولاً با مقدار بزرگی خطا مشخص میشود. برای برون رفت از این موضوع بایستی ماتریس وزن برای ویژگی‌ها رو طوری تغییر داد که در نهایت شبکه عصبی قادر به شناسایی دقیق اعداد باشد. این پرو سه از خطای ناشی از جواب فیدبک گرفته و بنا به این فیدبک ماتریس را دستخوش تغییر داده. این روش را نشر بازپس خطا^۲ گویند.

به طور کلی، پرسپترون یک مدل ساده از شبکه عصبی است. این در حالی است که مدل های چند لایه از شبکه عصبی بر اساس همین مدل ساده پرسپترون ساخته شده است. هنگامی که شبکه یادگیری لازم را برای شناسایی و طبقه بندی درست اشیاء را انجام داد، برنامه حاضر به اجرا خواهد بود و به صورت برنامه اجرایی (دیگر نیازی به یادگیری دوباره نیست) به اجرا خواهد آمد. از مثال های کارآمد متدهای DNN^۳ میتوان به شناسایی اعداد در دستخط ها برای استفاده از چک های بانکی در دستگاه های خودپرداز و یا شناسایی عکس های دوستان در سامانه امنیتی فیسبوک، پیشنهاد فیلم بر اساس 50

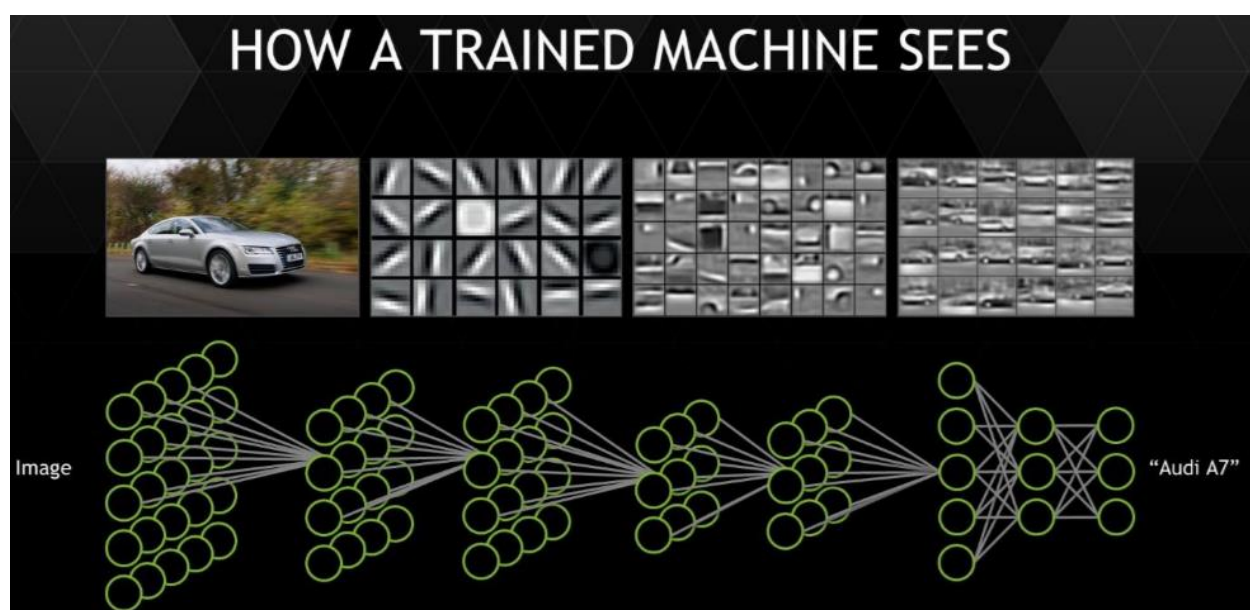
¹ Forward propagation

² Backward propagation

³ Deep neural network

میلیون کاربر شبکه اینترنتی نتفلیکس^۱، شناسایی و طبقه بندی انواع مختلف خودروها، پیاده ها و خطرهای جاده برای ماشین های بدون راننده و در نهایت ترجمه گفتار انسان به صورت بلادرنگ.

یک مدل چندلایه شبکه عصبی همانطور که در شکل زیر مشخص است، ممکن است از چند لایه که گره های آن به یکدیگر متصل هستند، تشکیل شده باشد؛ هر کدام از این گره ها مشخص کننده ی یک ویژگی ورودی است که خروجی هر کدام از این گره ها به عنوان ورودی چند گره دیگر در لایه های بعدی باشد.



شکل 27: یک مدل پیچیده از شبکه عصبی که نیازمند محاسبات سنگین است

در مدل شکل بالا، اولین لایه، از شبکه عصبی، ویژگی های اصلی از عکس را به قسمت های مختلف می شکند تا الگوهای اولیه از تصویر مانند خط ها و زوایا به دست آید. پس از آن در لایه دوم این با سر هم سازی این خطوط و زوایا الگویی پیچیده تر نسبت به لایه قبل مانند چرخ ها، شیشه ها و برف پاکن ها را می سازد و در لایه بعدی مدل خودرو را مشخص میکند و در نهایت در لایه آخر این شبکه عصبی برند خودرو را مشخص میکند. که در این مثال خودرو آوودی^۲ A7 بود.

¹Netflix

²Audi

یک راه دیگر به جای استفاده از شبکه کامل وصل^۱ استفاده از لایه های کانولوشنی^۲ است. یک عصب در یک شبکه کانولوشنی به یک عصب دیگر تنها از طریق یک ناحیه نسبتاً کوچکتر از مدل پیشگفته متصل خواهد بود. به طور کلی، این ناحیه ممکن است از یک شبکه ۵ در ۵ اعصاب تشکیل شده باشد (شاید هم ۷ در ۷ یا ۱۱ در ۱۱). اندازه این شبکه، اندازه فیلتر ما را مشخص میکند. بنابراین یک لایه کانولوشنی عملاً به عنوان یک تابع کانولوشن برای داده ورودی نقش ایفا میکند. این نوع از مدل شبکه کانولوشنی یک مدل ساده شده از مدل واقعی مغز واقع در شبکه گانگلونی^۳ است که عمدتاً وظیفه کورتکس دیداری را برعهده دارد.

در شبکه کانولوشنی DNN، ماتریس وزن فیلترها برای لایه های مختلف از شبکه عصبی مانند هم هستند. به دیگر سخن، یک شبکه کانولوشنی همانند یک مدل زیرلایه وار که هر کدام از فیلترهای مختلفی ساخته شده اند پیاده سازی می شوند صدها فیلتر مختلف ممکن است برای لایه کانولوشنی به کار روند. این متد ممکن است تا جایی پیش رود که یک لایه کانولوشنی از چندصد تابع کانولوشنی استفاده کند و نتیجه همه ی آنها را به لایه بعدی خود به عنوان ورودی دهد. شبکه های DNN که به صورت لایه های کانولوشنی به کار رفته اند را عموماً شبکه های عصبی کانولوشنی^۴ گویند.

4.2 پردازنده های گرافیکی NVidia: موتوری قوی برای یادگیری ژرف

با روی کار آمدن شبکه های عصبی ژرف و یا شبکه های کانولوشنی این سوال پیش آمد که محاسبات قوی و طولانی برای یادگیری این شبکه از طریق روش پس انتشار خطا چگونه میخواید انجام پذیرد. به علاوه، شبکه های عصبی ژرف نیازمند دیتاست های حجیم برای فرایند یادگیری هستند تا بتوانند در مرحله ی تست جوابی با دقت بالا در اختیار کاربر قرار دهند. این بدان معناست که هزار میلیون از نمونه های ورودی وجود دارد که بایستی هم از طریق روش انتشار رو به جلو و هم روش پس انتشار خطا مورد آزمایش قرار گیرند.

¹Fully connected

²Convolutional layer

³Retinal ganglion

⁴Convolutional Neural Networks (CNNs)

امروزه دیگر این امر مشخص است که یادگیری این شبکه ها چه در زمینه های صنعتی و چه در زمینه های آموزشی با استفاده از پردازنده های گرافیکی صورت میپذیرد، چرا که این پردازنده ها نسبت به مدل های هم تای خود یعنی پردازنده های اصلی (CPU) هم از لحاظ مصرف انرژی و هم از لحاظ سرعت اجرای برنامه ها به صورت موازی بسیار قوی تر و مقرون به صرفه تر هستند. برای روشن تر شدن این موضوع، لازم است گفته شود که شبکه های عصبی خود از مقداری زیادی عصب تشکیل شده اند و این اعصاب طبیعتاً به صورت موازی با هم در ارتباطند. این همکار موازی عصب ها با یکدیگر یک نگاشت ساده بر روی پردازنده های گرافیکی برای ما پدید می آورد که به مراتب سرعتی بالاتر در آموزش این شبکه ها نسبت به مدل سنتی یعنی استفاده از CPU در اختیار ما قرار میدهد.

شبکه های عصبی به صورت کلی بر اساس ماتریس هایی با ابعاد بزرگ و محاسبات سنگین در چند لایه ای سوار ده اند که محاسبات لازم برای این شبکه ها عصبی خود نیازمند پردازنده های قوی با قابلیت محاسبات ممیز اعشاری با پهنای باند های وسیع و سریع است. پردازنده های گرافیکی با تعدد هسته هایی محاسباتی ای که دارند، این محاسبات ماتریسی سنگین را برای ما بهینه کرده اند به طوری که کارایی در حدود چند صد TFLOPS در اختیار ما قرار داده اند.

در حال حاضر شرکت NVidia سردمدار پردازنده های گرافیکی است که به صورت اختصاصی معماری هایی برای شبکه های عصبی ژرف و هوش مصنوعی تولید میکند. پردازنده های گرافیکی شرکت NVidia به این شبکه های عصبی شتاب لازم را می دهند که برنامه های خود را با سرعت های ده تا بیست برابری به اجرا در آورند. این کار ها عموماً با کاهش زمان یادگیری برنامه از چندین هفته و چندین روز به تنها چندین ساعت است. با همکاری بسیاری از متخصصین در زمینه های مختلف، این امکان برای پردازنده های گرافیکی این شرکت به وجود آمده است که طراحی پردازنده های گرافیکی، معماری این سیستم ها، کامپایلر و الگوریتم های مورد نیاز را ارتقا ببخشند. تا سه سال قبل، پردازنده های گرافیکی شرکت NVidia کمک های شایانی را برای بهبود سرعت شبکه های یادگیری ژرف کرده اند که این رشد نمایی گویا قرار بر توقف به خود نخواهد گرفت.

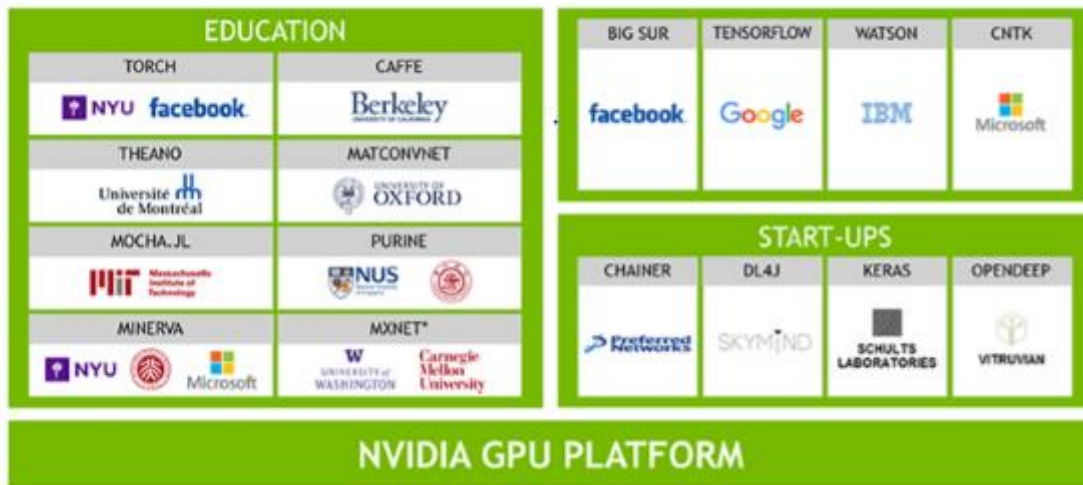
4.3 نرم افزار ها و toolkit های لازم برای یادگیری ژرف

همانطور که پیشتر گفته شد، هوش مصنوعی، یکی از سردم داران علم های به روز دنیاست. برنامه نویسی سخت در این زمینه برای برنامه نویسین و توسعه دهندگان این علم یکی از محدودیت های بزرگ است. برنامه نویسی نسبتا آسان و غنی پلتفرم کودا این امکان را به برنامه نویسان و محققان داده است که سرعت پیشرفت این علم را بیشتر کنند. شرکت NVidia برنامه ها و کتاب خوانه های مختلفی از قبیل NVIDIA DIGIT، cuDNN، cuBLAS و دیگر برنامه ها را برای اپلیکیشن های یادگیری ماشین به منظور اجرا بر روی بستر های ابری، دیتا سنترها، و workstation ها و یا حتی در سیستم های embedded و نهفته تحت عنوان Deep Learning software kit (SDK) فراهم کرده است. هدف اصلی برنامه نویسان آن است که بتوانند برنامه ای بنویسند که هر جایی که خواستند بتوانند آن را به اجرا در آورند. پردازنده های گرافیکی NVidia در نقطه از جهان وجود دارد. چه در سیستم های شخصی خانگی، چه در لپتاپ ها، سرورها و یا در سوپر کامپیوتر ها؛ و حتی در سیستم های مبتنی بر بستر ابری همچون آمازون¹، گوگل، IBM، فیسبوک، بایدو² و ماکروسافت. تمامی فریمورک های مبتنی بر هوش مصنوعی مبتنی بر پردازنده های گرافیکی NVidia هستند، از شرکت های اینترنتی گرفته تا استارتاپ های و موسسات تحقیقاتی. شرکت NVidia به منظور هر مصرفی از خانگی گرفته تا سیستم های سوپر کامپیوتر پردازنده گرافیکی خاص خود را تولید مینماید؛ برای مثال برای مصارف خانگی GeForce را و یا برای سیستم های ابری و یا سوپر کامپیوتر ها معماری Tesla را معرفی کرده است. برای ربات ها و یا پهباد ها نیز معماری Jetson که برای سیستم های نهفته کاربرد دارد را معرفی نموده است.

Amazon¹

Baidu²

ACCELERATE EVERY FRAMEWORK



*U. Washington, CMU, Stanford, Tsinghua, MIT, Microsoft, U. Alberta, MIT, NYU Shanghai

شکل 28: فریم وورک های مختلف برای کار با یادگیری ژرف

فصل چهارم

جمع‌بندی و نتیجه‌گیری و پیشنهادات

جمع‌بندی و نتیجه‌گیری

به طور کلی میتوان گفته که پردازنده های گرافیکی شرکت NVidia با رشد فزاینده ای که این سال ها در زمینه های گوناگون داشته اند، به بزرگترین پرچم دار سخت افزار های الگوریتم های یادگیری ژرف تبدیل شده اند. شرکت NVidia با تمرکزی خاص بر روی الگوریتم های یادگیری ژرف، هوش مصنوعی و شبکه های عصبی این امکان را به موسسات مختلف و دانشگاه ها داده است که فرایند یادگیری شبکه های عصبی را با سرعتی چند برابر حالت قبل به اجرا در آورند. گفتنی است که شرکت NVidia با ارائه نسل جدید از معماری پردازنده گرافیکی تحت عنوان پاسکال در را برای دیگر پردازش های موازی نیز بازتر کرد و این امکان را به کاربران داده است که کار با اعداد اعشاری را با سرعتی دو برابر نسل قبلی یعنی کپلر به اجرا در آورد. این شرکت، با ارائه تکنولوژی جدید NVLink از محدودیهایی که درگاه PCI-E خلاص شده و با سرعتی 4 برابر این درگاه، ارتباطات بین پردازنده های گرافیکی را برقرار ساخته است. از دیگر امکاناتی که این معماری همراه با پلتفرم کودا 8 ارائه کرده است، نسل جدیدی از حافظه های یکپارچه میبا شد که کاربر را از تخصیص دهی حافظه به صورت تخصصی رهایی میبخشد و عملاً این امکان را به حافظه اصلی و حافظه گرافیکی میدهد که با هم و در یک زمان به یک نقطه از حافظه دسترسی داشته باشند. گفتنی است که این حالت در معماری های قبلی اگر رخ میداد، باعث خطا در واکشی داده ها از حافظه شده و برنامه خطا نشان میداد.

پیشنهادهای

پردازنده های گرافیکی NVidia به منظور جلب رضایت برنامه نویسان نیازمند پلتفرمی کاربر دوستانه تر است، چرا که هر برنامه ای که به زبان کودا نوشته می شود تنها مختص به معماری ای است که برنامه نویس برای آن نوشته است، میباید. فلذا این پلتفرم باید این امکان را برای کاربران و برنامه نویسان فراهم کند که فارغ از معماری فعال در سیستم، در بهینه ترین حالت ممکن، برنامه به اجرا در آید. از دیگر پیشنهادهای که به این معماری گرافیکی میتواند وارد کرد، عدم اجرای برنامه هایی است که بر روی معماری های اصلی اجرا می شوند. این پردازنده هر چند به منظور کارهای گرافیکی طراحی شده است ولی در آینده ای نه چندان دور میتواند جایگزینی مناسب برای پردازنده های اصلی باشد.

منابع و مراجع

- [1] دیتاشیت معماری GP100 شرکت NVidia
- [2] دیتاشیت معماری jetron از شرکت NVidia.
- [3] Professional CUDA C Programming
Book by John Cheng, Max Grossman, and Ty McKercher
- [4] CUDA for Engineers: An Introduction to High-Performance Parallel
Computing
Duane Storti, Mete Yurtoglu