

Rendering Pipeline

مقدمه

زمانی که معماری پردازنده‌های گرافیکی تغییر کرد و به معماری Streaming تبدیل شد، شکل محاسبه و نحوه ارتباطات بر روی بردها تغییر زیادی یافت و پس از آن پردازنده‌های گرافیکی به انجام چند صد گیگافلاپ عملیات Single Precision Floating-Point (به شکلی از دقت نمایش اعداد اعشاری گفته می‌شود که اعداد توسط یک عبارت کامپیوتری قابل نمایش هستند) قادر شدند، این درحالی بود که قوی‌ترین پردازنده‌های آن زمان تنها در حدود 12 گیگافلاپ قدرت داشتند. البته این اختلاف هر روز زیادتر شده و می‌شود و امروز به موقعیتی رسیده‌ایم که کارت گرافیکی GTX 295 در حدود 2 هزار گیگافلاپ قدرت دارد این در حالی است که قدرت پردازنده‌ی گران‌قیمت شرکت اینتل به نام Intel Core i7 965 Extreme تنها در حدود 70 گیگافلاپ است. در این شماره قصد داریم بخشی از مقاله نوشته شده توسط Emmet Kilgariff و Randima Fernando از شرکت انویدیا در مورد معماری پردازنده‌های گرافیکی جی‌فورس سری شش را مورد بررسی قرار دهیم، در این مقاله به بررسی برترین محصول این سری یعنی Geforce 6800 خواهیم پرداخت. اگرچه این سری از پردازنده‌های گرافیکی کمی قدیمی هستند اما با مطالعه‌ی این مقاله به درک بهتری از معماری کارت‌های گرافیکی خواهید رسید.

نحوه‌ی قرارگیری یک پردازنده‌ی گرافیکی در سیستم کامپیوترها

در کامپیوترهای مدرن امروزی، پردازنده مرکزی با پردازنده گرافیکی از طریق اتصال PCI-Express یا AGP موجود بر روی مادربرد ارتباط برقرار می‌کند. از آنجایی که این اتصال مسئول انتقال داده‌های مربوط به فرمان‌ها، Textureها و Vertexها از پردازنده مرکزی به پردازنده گرافیکی است، این پهنای باند آن در سال‌های اخیر در کنار پردازنده گرافیکی پیشرفت کرده است. اسلات AGP، 66 مگاهرتز سرعت دارد و داده‌ها را با طول 32 بیت جا به جا می‌کند و در مجموع 264 MB/sec انتقال داده را فراهم می‌آورد. پس از AGP، نسل‌های بعدی این اسلات با نام‌های AGP 2x, 4x, 8x عرضه شدند که با ارائه هر کدام از این نسل‌ها پهنای باند اتصال دوبرابر شد تا این که در سال 2004

اسلاتی با حداکثر پهنای باند ممکن از لحاظ تئوری، با نام PCI-Express عرضه گشت. اسلاتی با پهنای باند، 4 GB/sec است.



شکل 1- یک پردازنده ی Geforce 6800

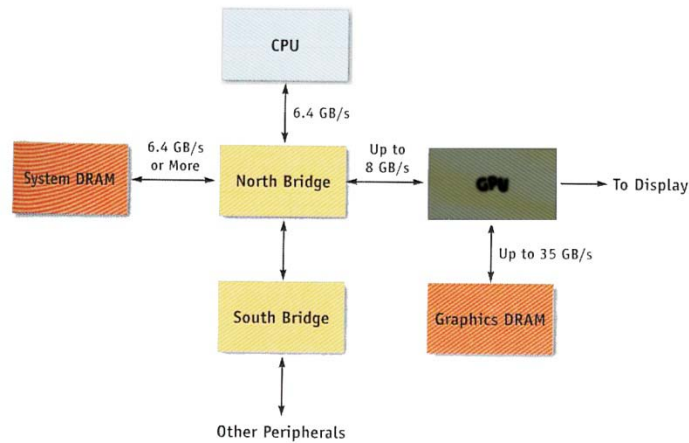
در جدول زیر قصد داریم تفاوت میان پهنای باند واسط حافظه پردازنده گرافیکی را با پهنای باند قسمت‌های دیگر سیستم مقایسه کنیم:

Component	Bandwidth(GB/sec)
واسط حافظه پردازنده گرافیکی	35 GB/sec
PCI Express Bus(x16)	8 GB/sec
واسط حافظه پردازنده مرکزی (800 Mhz FSB)	6.4 GB/sec

این جدول به شکل ویژه، تفاوت محسوس پهنای باند موجود در داخل پردازنده گرافیکی با پهنای باند قسمت‌های دیگر سیستم را نمایش می‌دهد و اهمیت الگوریتم‌هایی را که بر روی پردازنده گرافیکی اجرا می‌شوند و از این پهنای باند عظیم بهره می‌برند را نشان می‌دهد.

نحوه قرار گرفتن پردازنده گرافیکی در معماری سیستم

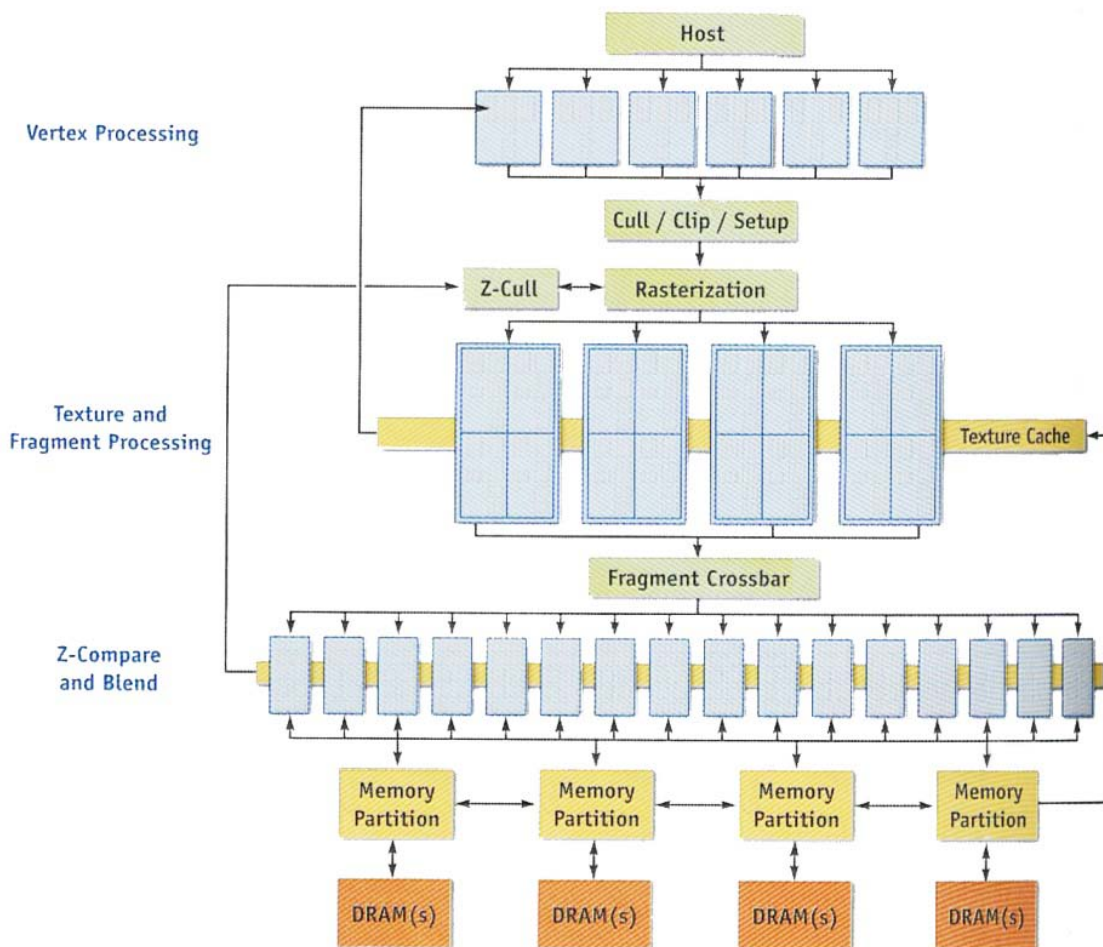
شکل 2 نحوه قرار گرفتن پردازنده گرافیکی در سیستم را نشان می‌دهد:



شکل 2- نحوه قرار گرفتن پردازنده گرافیکی در معماری سیستم

بلاک دیاگرام نشان دهنده عملیات‌های گرافیکی که انجام می‌شوند

شکل 3 معماری پردازنده گرافیکی Geforce 6 را نمایش می‌دهد. در این قسمت قصد داریم به بررسی Graphics Pipeline بپردازیم. این Pipeline ورودی‌هایی که توسط CPU ارسال می‌شود را دریافت کرده و پیکسل‌های خروجی را بر روی FrameBuffer نقش می‌بندد.



شکل 3- بلاک دیاگرام معماری Geforce6

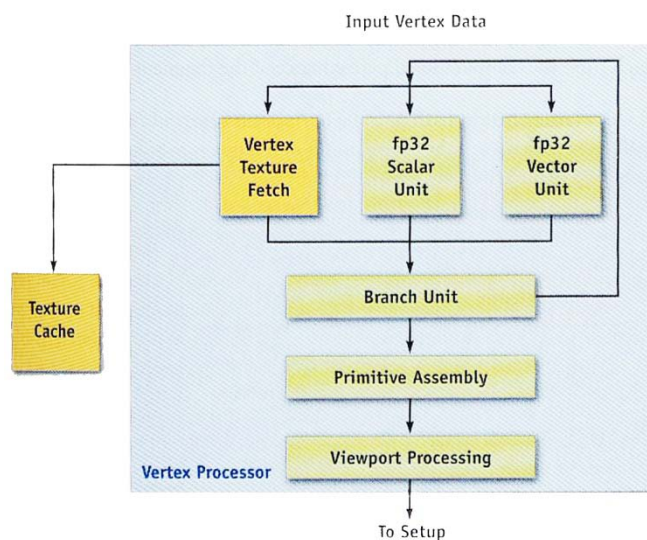
در ابتدا فرمان‌ها، Textureها و داده‌های مربوط به هر Vertex از طریق بافر مشترک بین پردازنده مرکزی و پردازنده گرافیکی موجود در حافظه سیستم یا حافظه مربوط به فریم‌بافر محلی، توسط پردازنده گرافیکی دریافت می‌شوند. زنجیره‌ای از دستورات توسط پردازنده مرکزی در این حافظه مشترک نوشته می‌شود که وظیفه مقداردهی اولیه و ویرایش state مربوط به عملیات‌های گرافیکی، ارسال فرمان‌های لازم جهت رندر کردن و ارجاع‌هایی به داده‌های مربوط به Vertexها و Textureها را دارد. فرمان‌ها به اجزا کوچکتر تقسیم می‌شوند و از واحد Vertex Fetch برای خواندن و ذخیره کردن اطلاعات مربوط به Vertexهایی که در فرمان‌های مربوط به عملیات رندر، به آنها ارجاع داده شده است، استفاده می‌شود. فرمان‌ها، Vertexها و Stateها هنگام عبور از Pipeline Stageهای متوالی، تغییر می‌کنند. پردازنده

مخصوصاً Vertex ها که Vertex Shader نامیده می‌شود (شکل 4)، این امکان را فراهم می‌کند که یک برنامه بر روی هر کدام از Vertex های مربوط به یک شی اجرا شود و تبدیلات، عملیات Skinning و دیگر عملیات‌هایی که بر روی Vertex ها انجام می‌شوند را اعمال کند. برای اولین بار در کارت گرافیک‌های سری شش شرکت انویدیا، این امکان فراهم شده است تا برنامه‌هایی که بر روی Vertex ها اجرا می‌شوند، بتوانند داده‌های مربوط به Texture ها را واکشی (Fetch) کنند. همه عملیات‌ها به صورت عملیات‌های 32 بیت، Floating-Point (fp32) موجود برای هر جز انجام می‌شوند. تعداد پردازنده‌های Vertex که همراه این معماری استفاده می‌شوند قابل تغییر است و مدل‌های دارای کارایی بیشتر و گرانتر از 6 واحد پردازش Vertex بهره می‌برند و مدل‌های ارزان قیمتتر از 2 واحد استفاده می‌کنند.

از آنجایی که پردازنده‌های Vertex، می‌توانند به Texture ها دسترسی پیدا کنند، موتورهای Vertex به حافظه‌ی کش مربوط به Texture ها متصل هستند. همچنین حافظه‌ی کش دیگری وجود دارد که داده‌های مربوط به Vertex ها را قبل و بعد از انجام پردازش توسط پردازنده‌های Vertex ذخیره می‌کند تا تعداد عملیات واکشی و محاسبات کاهش یابند، به این ترتیب که اگر اندیس یک Vertex در استفاده از تابع Draw دوبار تکرار شود (مثلاً در هنگام استفاده از روش اتصال Strip، برای ساخت یک سه ضلعی)، برنامه‌ای که بر روی Vertex عمل می‌کند، نیازی به انجام دوباره عملیات بر روی Vertex نخواهد داشت و از نتیجه موجود در حافظه‌ی کش استفاده می‌کند.

سپس Vertex ها گروه‌بندی شده و باعث شکل‌گیری اشکال ابتدایی نقطه، خط و سه ضلعی می‌شوند. سپس بلاک‌های Cull، Clip و Setup، عملیات‌هایی را بر روی هر شکل ابتدایی انجام می‌دهند. بلاکی که مسئولیت عملیات Culling را برعهده دارد، اشکال ابتدایی قابل رویت نیستند را حذف می‌کند، بلاک Clipping، اشکال ابتدایی که با هر ناقص شکل مشاهده صحنه، برخورد می‌کنند را مشخص می‌نماید و بلاک Setup، معادلات مربوط به لبه‌ها و صفحه‌ها را به گونه‌ای سازماندهی می‌کند که داده‌ها برای عملیات Rasterization آماده شوند. بلاک Rasterization محاسبه می‌کند کدام پیکسل و یا Sample (اگر MultiSampling فعال باشد) توسط کدام شکل ابتدایی

پوشانده شده است و به سرعت با استفاده از بلاک z-cull برای دور انداختن پیکسل‌ها و یا Sample‌هایی که توسط شی دیگری که دارای عمق کمتری است، پوشانده شده‌اند، اقدام می‌کند. پیکسلی که عملیات بر روی آن انجام خواهد شد به عنوان Fragment شناخته می‌شود. این Fragment، از Fragment Processor و تست‌های مختلف عبور می‌کند و اگر از تمام آن‌ها با موفقیت خارج شود، حاوی رنگ و عمق مربوط به پیکسل موجود بر روی فریم‌بافر و یا Render Target خواهد بود.



شکل 4- پردازنده Vertex مربوط به Geforce6

شکل 5 نشان دهنده یک Fragment Processor (که معمولاً با نام Pixel Shader شناخته می‌شود) و یک Texture Pipeline است. واحدهای Fragment Processing و Texture Processing، در کنار یکدیگر عمل می‌کنند تا یک برنامه shader را بر روی هر Fragment به صورت مجزا اجرا کنند. معماری Geforce6 به گونه‌ای است که می‌توان به آن تعداد متغیری Fragment Processor متصل کرد یعنی پردازنده‌های گرافیکی سری Geforce 6 دارای تعداد متفاوتی Fragment Pipeline یا Pixel Pipeline می‌باشند. همانند Vertex Processorها، داده‌های مربوط به Textureها نیز بر روی تراشه، در درون کش ذخیره می‌شوند تا از پهنای باند خارجی کمتری استفاده نمایند و کارایی افزایش یابد.

واحدهای پردازش Texture و Fragment در هر لحظه بر روی چهار ضلعی‌هایی که از چهار پیکسل تشکیل شده‌اند (که Quad نامیده می‌شوند) عمل کرده و مشتقاتی از آن را محاسبه می‌کنند که برای تعیین سطح

جزئیات مربوط به آن بخش از Texture، قابل استفاده است. همچنین Fragment Processor در هر لحظه بر روی گروه‌هایی متشکل از صدها پیکسل عمل می‌کند و با یک دستور همزمان پردازش را بر روی چندین داده انجام می‌دهد به این پردازنده‌ها SIMD (Single Instruction Multiple Data) گفته می‌شود (که در آن یک موتور پردازش Fragment در هر لحظه فقط بر روی یک Fragment کار می‌کند) و با این کار تاخیر مربوط به عملیات واکشی، Texture برکارآیی محاسباتی پردازنده Fragment تاثیری نخواهد داشت.

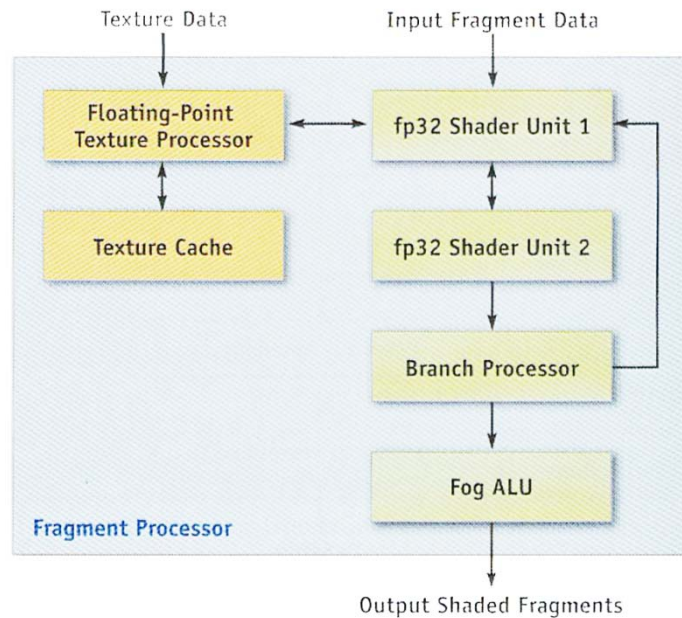
پردازنده Fragment، از واحد Texture برای واکشی داده از حافظه استفاده می‌کند. واحد Texture از فرمت‌های مختلفی پشتیبانی می‌کند (فرمت تعیین کننده شکل ذخیره شدن جزئیات هر پیکسل یا نحوه کاربرد آن می‌باشد)، سپس این داده‌ها می‌توانند توسط روش فیلترینگ bilinear، trilinear یا Anisotropic، فیلتر شوند. همه داده‌ها در نهایت به شکل fp16 یا fp32 به پردازنده Fragment ارسال می‌شوند. یک Texture را می‌توان به صورت آرایه‌ای دوبعدی یا سه بعدی از داده‌ها که توسط یک واحد Texture در محل‌های دلخواه حافظه قرار می‌گیرند فرض کرد و تحت عملیات فیلترینگ به Texture دیگری تبدیل کرد که مجدداً می‌تواند توسط همین بخش عملیاتی به Texture دیگری تبدیل شود. این بخش عملیاتی به صورت متوالی می‌تواند بر روی یک Texture عمل کند. سری Geforce6 قادر است عملیات فیلترینگ را بر روی Texture‌هایی با فرمت fp16 به صورت سخت‌افزاری انجام دهد.

واحد پردازش Fragment دارای دو واحد fp32 shader در هر Pipeline است و Fragment‌ها قبل از این که مجدداً به Pipeline برگردند از این دو واحد fp32 shader و Branch Processor عبور می‌کنند و سپس مجدداً به Pipeline برگشته تا سری بعدی دستورات بر روی آن‌ها اعمال شود. این عملیات بازگشت در هر کلاک فقط یک بار انجام می‌شود. همچنین می‌توان از واحد fp32 shader اول، برای اصلاح مختصات Texture مربوط به هر vertex جهت ایجاد پرسپکتیو در صحنه استفاده نمود یا برای انجام عملیات‌های ضرب از آن بهره برد. می‌توان هشت یا تعداد بیشتری عملیات ریاضی را در یک کلاک با Pixel Shader انجام داد و اگر عملیات واکشی texture در واحد اول shader رخ دهد، می‌توان 4 عملیات ریاضی را در هر کلاک انجام داد.

در آخرین مرحله از Pixel Shader Pipeline، واحد ایجاد مه برای ترکیب مه و تصویر استفاده می‌شود، بدون این که کارآیی پردازنده گرافیکی، ذره‌ای کاهش پیدا کند. عملیات ترکیب مه، و تصویر حاصل شده تا این مرحله، در نرم افزارهای گرافیکی به وسیله تابع زیر انجام می‌شود:

$$\text{رنگ مه} \times \text{میزان رنگ مه} + \text{رنگ پیکسل حاصل شده تا این مرحله} \times (\text{میزان رنگ مه} - 1) = \text{رنگ پیکسل خروجی}$$

Fragment هایی که از واحد پردازش Fragment خارج می‌شوند به همان ترتیبی که Rasterize شده‌اند به واحد مقایسه z و Blend فرستاده می‌شوند که در آنجا عملیات تست عمق (مقایسه مقدار z و به روزرسانی آن)، عملیات‌های Stencil، ترکیب آلفا و نوشتن رنگ نهایی در صفحه مقصد (یک render target یا بافر رندر که به صورت Off-Screen است و بعد از این که کاملاً رندر شد on-Screen می‌شود- بر روی صفحه نمایش به تصویر درمی‌آید-) می‌باشد انجام می‌شود.



شکل 5- پردازنده‌ی Fragment و Texel Pipeline مربوط به Geforce6

حافظه سیستم حداکثر به چهار قسمت تقسیم شده است و هر کدام از آن‌ها به DRAM هایی متصل هستند. پردازنده‌های گرافیکی به جای استفاده از RAM های متعارف در بازار از DRAM های استاندارد بهره می‌برند. این کار باعث صرفه‌جویی و کاهش

قیمت کارت گرافیکی می‌شود. استفاده از قطعات کوچکتر و مستقل حافظه، باعث می‌شود تا سیستم حافظه در هر صورت چه در حالت انتقال بلاک‌های بزرگ و چه بلاک‌های کوچک از حافظه بهینه‌تر عمل کند. همه صفحه‌های رندر شده در DRAM ذخیره شده اند در حالی که داده‌های مربوط به Textureها و ورودی‌ها می‌توانند در DRAMهای پردازنده گرافیکی یا سیستم ذخیره شوند. حافظه چهار تکه، برای پردازنده گرافیکی پهنای باندی معادل 256-bit فراهم می‌کند. سیستم مربوط به این حافظه، امکان ارسال داده‌های نسبتاً کوچک 32 بیتی به صورت زنجیره‌ای با محدودیت سخت‌افزاری 35 GB/sec را فراهم می‌کند.