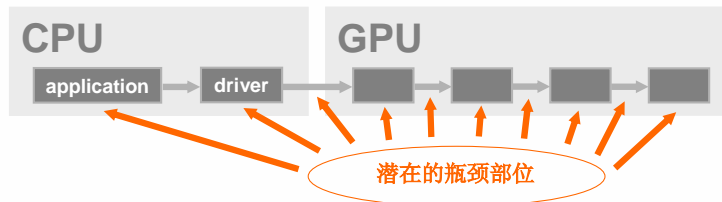




Hi. My name is Koji Ashida. I'm part of the developer technology group at NVIDIA and we are the guys who tend to assist game developers in creating new effects, optimizing their graphics engine and taking care of bugs. Today, I'm going to talk to you about optimizing the graphics pipeline.

概述



- 瓶颈（bottleneck）制约着总的吞吐率
- 总的来说，每个应用程序甚至每一帧都会有不同的瓶颈部位
- 对于流水线（pipeline）架构来说，要获得好的性能取决于发现和排除瓶颈的影响



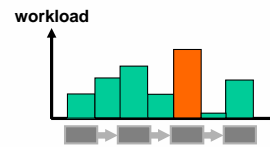
#VIDIA.

We're going to go through all the stages of the GPU plus talk a little bit about what happens on the CPU side that can bottleneck a graphics application. As we all know, in order to render a scene, the application which is running on the CPU must send data and instructions across to the device. It communicates with the device through the driver. Then, once the device has the data, it processes the data using the graphics chip itself, and finally writes it out to the frame buffer. Because this whole process is a single pipeline from the CPU to the last stage of the GPU, any one of those stages can be a potential bottleneck. The good thing about a pipeline is that it's very efficient and, in particular, it's very efficient at rendering graphics because it can parallelize a lot of operations. The bad thing about a pipeline is that, once you do have a bottleneck, then the whole pipeline is running at that speed so you really want to basically level off all your stages in the pipeline such that they have equal workloads. Hopefully, by the end of this presentation, you'll have a very good idea of how to either reduce the workload on a certain stage or at least increase the workload on the other stage such that you're getting better visual quality.

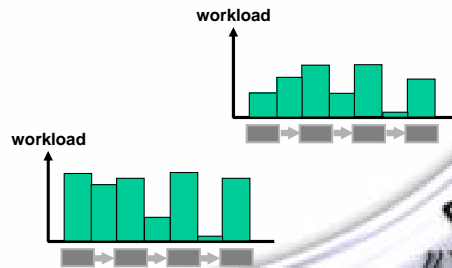
瓶颈的定位和排除

Locating and eliminating bottlenecks

- 定位：在每一个阶段（stage）
 - 改变它的负荷
 - 总的性能受到了明显的影响吗？
 - 降低时钟
 - 总的性能受到了明显的影响吗？



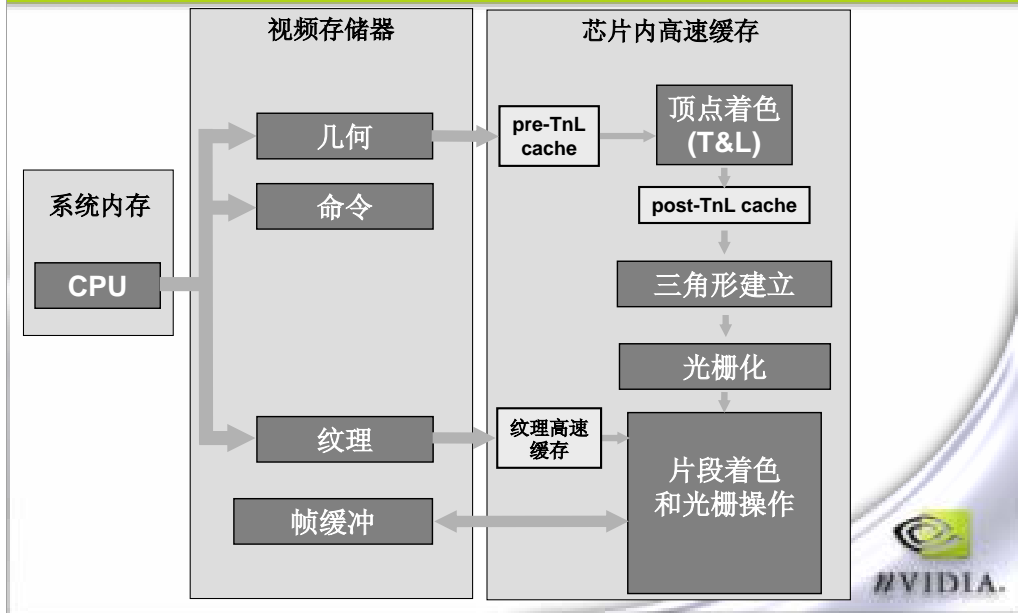
- 排除：
 - 降低产生瓶颈的部位的负荷
 - 增加未发生瓶颈部位的负荷



#VIDIA.

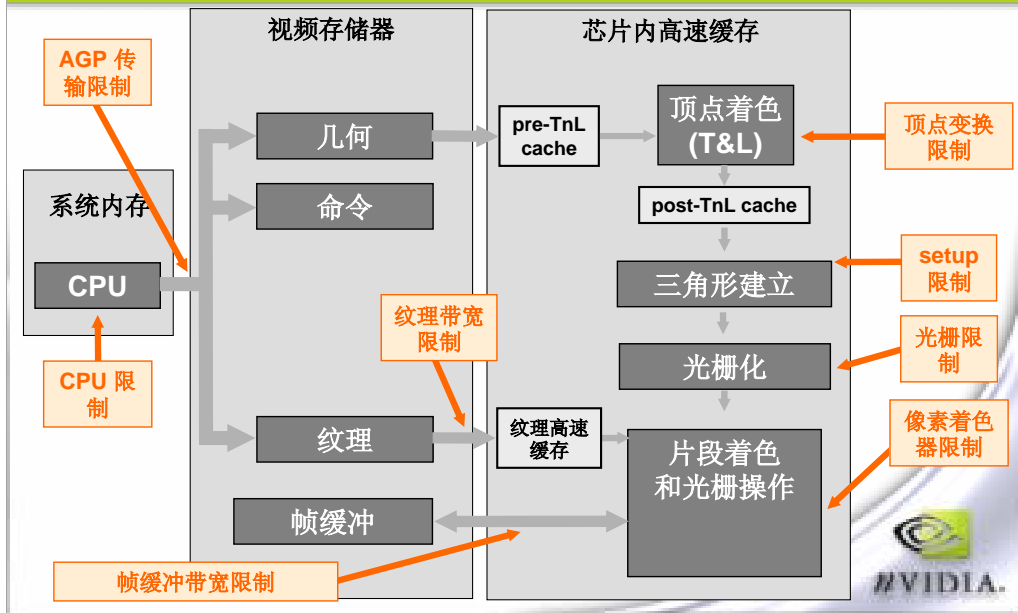
Now, please remember that, for any given scene, you're going to have different bottlenecks for different objects, different materials, different parts of the scene and so it can be fairly difficult to find where the impact is. Now, you do have two choices once you've identified the bottleneck. One thing you can do is try to reduce the overall workload of that stage and thereby increasing the frame rates or what you can do is increase the workload of all the other stages and, therefore, increasing the visual quality. Now, you want to choose which one you do, depending on your target frame rate so, if the application is already running at 80 frames per second, you may, instead of seeking to run at 100 frames per second, increase the workload on the other stages. The basic theory that we're going to use is that we're going to step through each stage and vary its workload. If it is the bottleneck, then the overall frame rate's going to change. If it isn't, then we're going to see no difference, supposedly.

图形渲染流水线 Graphics rendering pipeline



This is an overall view of the graphics pipeline. To the far left, we see the CPU, which is where the application and driver are going to be running, and the CPU is communicating with the graphics device through the AGP bus. These days, we have AGP8X, so it's running pretty fast. It communicates both with the graphics chip itself and with video memory through the graphics chip memory controller. In video memory, typically what's stored is static geometry or semi-static geometry, also the command stream, textures, preferably compressed, and, of course, your frame buffer and any other intermediate services that you have. Then, on the actual hardware chip, we have some caches to do buffering and to ensure that the pipeline's running as optimally as possible.

潜在的瓶颈部位 Potential bottlenecks



The first real module that we encounter is the vertex Shader, also called the vertex program units, and this is where you're going to do your transform and lighting of vertices. These vertices then get put into some sort of vertex cache that operates in different fashions. Some operate as FIFO's, some, as leased recently used, or LRU, and, obviously, current high-end cards have larger caches than the older generation and mainstream cards. Then the triangle setup stage is reached. This module reads vertices from the cache and this is where basically the polygon is formed. Once the polygon is formed, the rasterization is where it gets broken up into pixels. So now the next stage only accepts pixels and this is the fragment shading or pixel shader stage. The pixel shader stage is typically where a lot of time is spent these days and we're going to spend a good amount of time analyzing this. After the pixel shading stage is the raster operation meaning the alpha blend and stencil Z buffer. These can contribute to the bottleneck but typically not. And, finally, you can have traffic from the pixel Shader and raster module to the frame by frame back, right, and back because you can do alpha blending, you need to read stenciling to read a Z and also the fragment Shader can read textures so you can have a texture bottleneck as you're accessing memory again.

图形渲染流水线的瓶颈

Graphics rendering pipeline bottlenecks

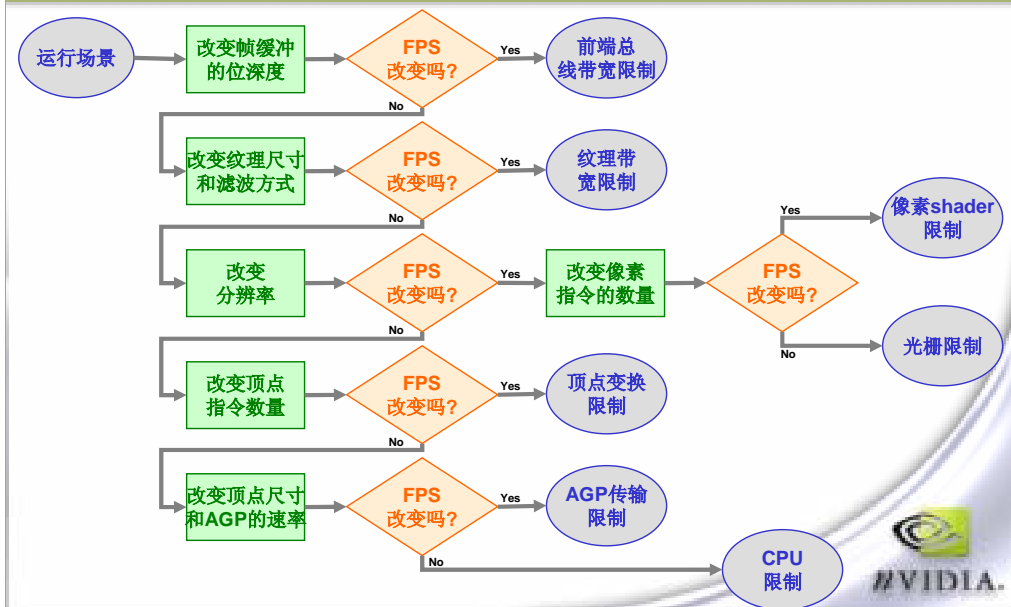
- 术语“**转换限制**（**transform bound**）”常常意味着瓶颈出现在光栅化之前的任何阶段
- 术语“**填充限制**（**fill bound**）”常常意味着瓶颈出现在“**setup**后的任何阶段”
- 转换和填充限制可以在一帧处理过程中同时存在！



#VIDIA.

When we're talking about bottlenecks, we tend to oversimplify the pipeline. We talk about being transfer-bound or fill-bound and, when we mean transfer-bound, it's really-- we're saying that it's anything before the polygon gets broken up into pixels; i.e., before they're rasterized. And when we mean fill-bound, we mean a number of things that happen after the polygon is formed, and that could be raster, texture-bound, fragment Shader bound or frame buffer bound.

瓶颈辨识 Bottleneck identification

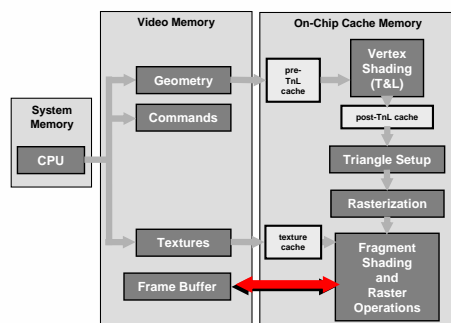


This is a chart that you can use to identify what part of the graphics pipeline is your bottleneck so I put the slide into the presentation just so you can have it as reference and then, over the next few slides, I'm going to explain each block.

帧缓冲带宽限制

Frame buffer bandwidth limited

- 改变所有渲染对象的颜色深度 (16-bit vs. 32-bit)
 - 如果帧率改变了，应用程序就是帧缓冲带宽限制的



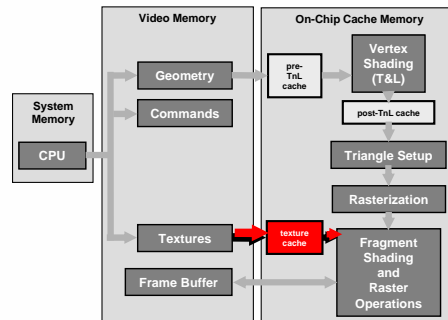
#VIDIA.

So let's start at the bottom of the pipeline. This is the part where the graphics engine is trying to read from the frame or to frame. Simply to access off-screen surfaces, not necessarily textures. So the easiest way to identify if this is your bottleneck is to vary your bit depth. So if you're running at 32 bits per pixel, run it at 16. If the frame rate varies then you're probably frame buffer limited. This isn't necessarily the case all the time but it does happen.

纹理带宽限制

Texture bandwidth limited

- 否则，改变纹理大小或者是纹理滤波类型
- 强制MIPMAP LOD（细节等级）偏量到+10
 - Point filtering vs. bilinear vs. trilinear vs. anisotropic
 - 如果帧率改变了，应用程序就是纹理带宽限制的

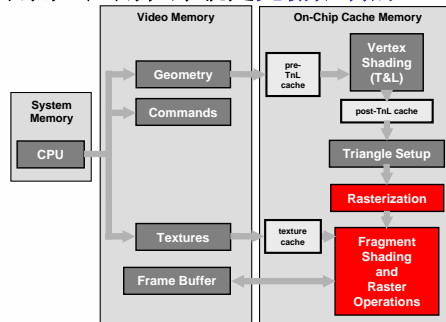


#VIDIA.

Let's move a little further up and we hit the fragment Shader but, even more importantly, the data that the fragment shader is pulling. This data is texels and texels can be a bottleneck if there are too many of them or if they don't hit the texture cache efficiently or if they're being filtered with an expensive filter that may not be improving the visual quality of the scene. So if you vary the LOD of your MIP maps and basically force them all to the smallest size possible, you'll be able to figure out if the texture resolutions are not adequate and I do hope that you are using MIP maps because it's very important in today's hardware. Next, you can vary the filtering. So if you're using bilinear or tri-linear, you can drop down to point filtering, for example, or you can increase the anisotropic ratio to 8X and, again, just check the frame rate.

片段或者光栅限制 Fragment or raster limited

- 否则，改变所有渲染对象的分辨率
- 如果帧率改变了，改变片段（**fragments**，或者说待处理像素）程序（**fragment programs**）的指令数量
 - 如果帧率改变了，应用程序就是受限于片段着色器（**fragment shader**）
 - 否则，应用程序就是光栅限制的



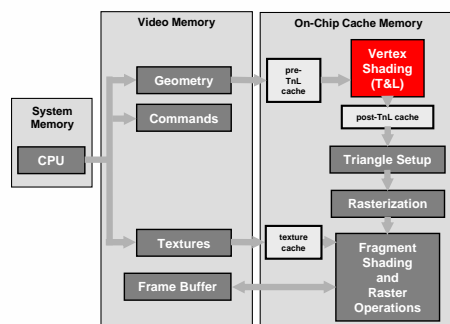
#VIDIA.

Next, we have the fragment program and the rasterization module. The first thing you should do here is vary your render target size, not the bit depth but the actual resolution so it'll go from 640x480 to 1024 to 1280 to 1600, for example. And if your frame rate varies, then you know that your bottleneck is somewhere here, somewhere in the fragment shader or the rasterizer. And to check if you're limited by the fragment shader, you can simply vary the number of fragment shader instructions. If the frame rate doesn't vary, then you are raster limited.

顶点变换限制

Vertex transform limited

- 否则，改变您的顶点程序（**vertex programs**）的指令数量
 - 注意：不要增加那些很容易被优化掉的指令
 - 如果帧率改变了，应用程序就是**顶点转换限制（vertex transform limited）**的



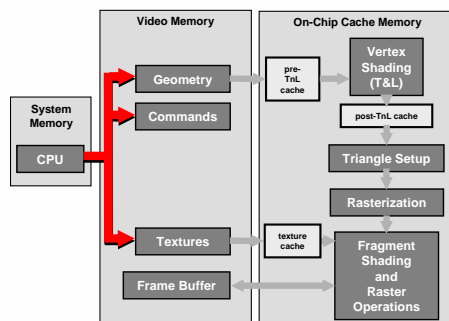
#VIDIA.

Next up the pipeline, we have the vertex engine and the quickest way to vary the load on that is to add or remove the vertex shader instructions. The only caveat here is that the driver compiler and optimizer can easily detect if you are doing operations that are optimizable and therefore easily removed. So try to have dependent instructions and write out some random color to the diffuse register. Again, if the frame rate varies, then this is where your bottleneck is.

AGP传输限制

AGP transfer limited

- 否则，改变顶点格式尺寸或者AGP传输率
 - 如果帧率改变了，应用程序就是AGP传输限制的

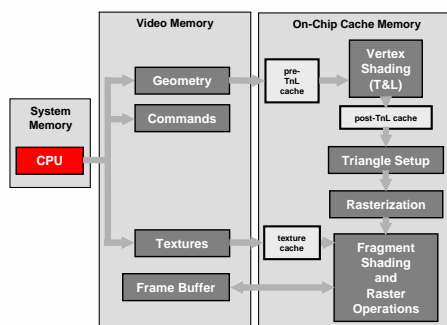


#VIDIA.

Even further up the pipeline, we have our AGP bus and this used to be a large bottleneck. If you have a lot of dynamic geometry, this could still be a bottleneck for you and, in that case, you can test it out by varying your vertex format size. Typically, you want to keep vertices to 32 byte multiples, so 32, 64, et cetera.

CPU限制 CPU limited

- 否则，应用程序受限于CPU



#VIDIA.

Finally, if none of that varied your frame rate, then you are CPU limited.

瓶颈鉴别捷径

Bottleneck identification shortcuts

- 在不同的**CPU**上使用同一个**GPU**运行
 - 如果帧率改变了，应用程序就是**CPU限制的**
 - 如果帧率和**CPU**的速度成正比则完全是**CPU限制**
- 在**BIOS**中强制让**AGP**运行在**1X**模式
 - 如果帧率改变了，应用程序就是**AGP带宽限制的**
- **降低**您的**GPU**的运行频率
 - 如果较低的运行频率影响了性能，应用程序就是**顶点变换、光栅或者片段着色器限制的**
 - 如果较慢的显示存储器的速度影响了性能，则应用程序是**纹理或者帧缓冲带宽限制的**



#VIDIA.

Now, there are a couple of quick tricks you can use to figure out your bottlenecks. First thing to do is to run the hardware on different CPUs, and that way you're varying your CPU and you can immediately tell if it's CPU limited if, and only if, the percentage increases are the same; i.e., the megahertz percentage matches the performance percentage increase. Now, the other thing you can do for AGP bottleneck is to vary the AGP transfer rates in your BIOS so you can drop it down to 1X and see if that impacts your frame rate. The other thing is using Power Strip or some utility, you can under clock your GPU and it will obviously affect the whole pipeline so at least you'll be able to tell if it's somewhere within the graphics chip. You can also slow down the memory clocks and, in that case, you'll be slowing down the texture reads or frame buffer read/write.

瓶颈辨别功能

Bottleneck identification functionality

- 编写尽可能简单的应用程序
- 有针对性
 - 对象: 地形(terrain), 角色(character), 建筑物(building), vehicles
 - 渲染流程(Pass): 每个pass就有不同的瓶颈
 - 人工智能/物理现象仿真循环(AI/physics loop): 去除CPU过载
- 能自由地冻结在某帧
 - 提供一致性



#VIDIA.

Now, we've had trouble in the past identifying bottlenecks and what helps us a lot is if the application is written in such a way that you can toggle different parts of the scenes. For example, if you can remove the terrain or remove all static objects and just leave character onscreen, then you can actually profile just the characters or, if you remove everything else and just render terrain, you can profile just the terrain and what this is doing is trying to identify the different bottlenecks that are going to occur on your different materials and objects, because each one of those can have a different workload at different stages. The other thing that would be interesting is for each of those objects, if you had multi-pass rendering, toggling the different passes because each pass can also have a different bottleneck. And, finally, something that would be interesting to do is to remove the CPU overhead or the application use of the CPU by, for example, pausing the AI and physics loops. Doing this will also allow you to basically pause the game such that you have no other variables affecting your frame rate and, that way, have a steady scene with which you can toggle the different objects and identify the different bottlenecks for each of them.

总体优化：批处理（Batching）

I

- 消除小批作业（**small batches**）：
 - 在每个顶点缓冲区/阵列中使用成千个顶点
 - 每一次调用绘制尽可能多的三角形
 - 每次调用绘制成千个三角形
 - 大约**100k DIP/秒**就可以使得**3.0GHz Pentium4**达到饱和
 - **50fps**意味着每帧有**2K DIP**!
 - 这时**2k三角形/帧**或者**2M三角形/帧**的性能是一样的



#VIDIA.

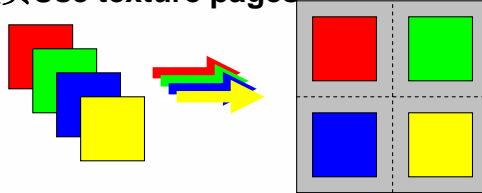
Now, there is one bottleneck on the CPU which is rather large and this has turned out to be batching. What we mean by batching is the number of draw primitive calls that you make. So a draw primitive call consists of sending a set of polygons to the driver, hence to the device, with a certain state, right, the state that you set before -- could be texture states, render states, or sampling states, et cetera. So that's what we define as a batch. Now, it turns out that, a P-4, 3 GHz CPU can only send 100,000 batches or 100,000 draw index primitive calls per second. What this means is that, at 50 frames per second, you only have 2,000 draw calls per frame. What this also means is that you can either choose to render 2,000 triangles per frame or 2 million triangles per frame. If you're using all of the 2,000 draw calls per frame, then ensure that you're sending a lot of triangles for every draw call.

总体优化：批处理（Batching）

II

消除小批作业 (继续):

- 使用退化的（**degenerate**）的三角形将条带（**strips**）连接在一起
 - 硬件剔除面积为零的三角形是**非常快速**的
- 使用纹理页**Use texture pages**



- 使用顶点着色器（**shader**）将需要的几何打包在一起。
 - VS2.0**和**VP30**具有**256个4D向量的常数**



#VIDIA.

Here are a couple of things that can help out with batching. You can join up objects or parts of objects by using degenerate triangles. Degenerate triangles are basically zero area triangles where two of the vertices are the same. Now, the culling engine in the chip removes these zero area triangles at basically no cost so don't be worried about sending inefficient triangles to the hardware. We'll take care of it. Next thing you can do is to use texture pages. You often render geometry with different draw calls because they use different materials and the materials use different textures. Well, a good trick that you could use is to store all your textures into a larger texture or effectively a texture page and the only thing you have to worry about then is the border around the texture because of MIP mapping and also filtering so the limitation is that you can't tile these textures so, for any other geometry that doesn't use tiling, this would be a good optimization. Lastly, it's a good idea to use vertex shaders to batch instanced geometry. What this means is that you can render several objects in one call, let's take a robot, for example. Each jointed segment of the robot is rendered with a different matrix. Usually we have to issue a draw call for each segment. What you can do, however, is that for each vertex, you send an index number which tells what matrix the vertex needs to be transformed with. And in the vertex shader, you find the correct matrix from the constant vector using the index, and use that to transform the vertex. This lets us to send all the vertices of a robot in one draw primitive call.

总体优化：索引和排序

Overall optimization: Indexing, sorting

- 使用索引的原始数据（条带化或者列表）
 - 这是唯一能够使用TnL前和后级Cache的唯一办法！
 - （非索引的条带化同样可以使用Cache）
- 重新对顶点排序以便有序地使用它们
 - 以获得最大化的cache使用率
- 将对象按照由外到里进行大致排序
- 按每个纹理以及渲染阶段对批进行排序



#NVIDIA

Another general rule for optimizing the engine is to use indexed primitives. Typically, strips are not long enough to efficiently render geometry all the time but indexed lists actually allow us to make best use of the vertex cache because non-indexed triangle lists do not use the vertex cache. The other thing that you want to do is to reorder your vertices such that they hit the same area of memory. Our chip has a DMA cache, and this cache stores vertices that were fetched from memory and if you're loading vertices from very different areas in memory, then you're going to be thrashing the cache and transferring more data than necessary. Next you want to order your scene front to back. What this is doing is taking advantage of our Z-cull technology. Z-cull has been around since GeForce 3 days and is also available in other people's hardware. What it does is accelerated removal of pixels that will never be seen and it's important to take advantage of this to get better performance, especially for scenes with high depth complexity. This can be actually used for a Z first technique that I'll talk about later. And, lastly, you want to sort your batches per texture because texture changes are the hardest hitter of all the state changes possible.

总体优化：闭塞面询问

Overall optimization: Occlusion query

- 使用闭塞面询问以保证顶点和像素的吞吐率：
 - 多pass的渲染：
 - 在第一个pass中就为每个对象附上询问（**query**）
 - 如果一个对象没有足够的像素被绘制，则直接忽略后面的
 - 粗略的可见度判断：
 - 画一个有查询能力的**quad**来获知太阳照射下的镜头眩光多寡
 - 绘制一个包含询问的**bounding box**以了解一个入口或者一个复杂的对象是否可见，如果不可见，则忽略掉对它的渲染。



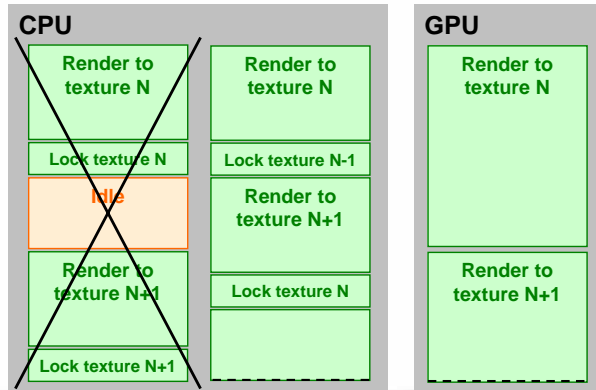
#NVIDIA

Recently exposed in the DX9 API (but it's been available in OpenGL for a while) are occlusion queries. Occlusion queries use an ID-based approach to query if an object tagged with this ID has been rendered. You let the hardware process it, meaning that there is no CPU intervention, but the GPU doesn't actually write anything, it just does a Z check and then it returns a flag to the driver that says, yes, something's been rendered, and gives you the number of pixels, or no, nothing's been rendered. So what you can do with this is try to render a very low-cost bounding box for your object after laying down large scene occluders, render the rest of the scene, then go back and query the driver to see if that object could have rendered any pixels. If it did, then send the real object down that will actually be visible. A way to optimize this is, in the case of multi-pass rendering, just send down the first pass but tag it. And, at the end of the scene, draw the subsequent passes if, and only if, it passed the occlusion test. The other thing you can use this for is something like a sun glare, where you don't want to be locking the frame buffer to check if the sun pixels were drawn or not. What you can do is render the sun with an occlusion tag and then, afterwards, render the glares, depending if the sun was visible or not and, because we returned the number of pixels drawn, you can basically fade in and fade out the sun.

总体优化：小心资源锁定！

Overall optimization: Beware of resource locking!

- 一个锁定资源的调用(**Lock, glReadPixels**)如果放错了位置的话会成为潜在的一个阻塞部位：
 - CPU处于空闲，等待着GPU 以继续处理
- 如果可能的话要消除这种情况
- 否则就放置它们因此GPU就有时间继续处理：



Resource locking basically means that the GPU is waiting to access a resource that, for some reason, you've grabbed. You've told us that you want to read from it, for example. This is bad. This will start synchronizing the GPU to the CPU. You want to avoid this at all cost. However, if you can't live with that, then what you can do is locking a resource that will not be used for another few frames and, that way, you ensure that the graphics unit is not waiting around for you to unlock it.

CPU瓶颈：原因

CPU bottlenecks: Causes

- 应用程序限制：
 - 游戏逻辑、AI、网络、文件I/O
 - 图形将受限于简单拣选和分类
- 驱动程序或者API限制：肯定有地方出了问题！
 - 没有采用快速途径（fast path）
 - 不正确地使用API
 - 有小批处理
- 大多数图形应用程序都是CPU限制型的



#NVIDIA

Now, let's move on to potential CPU bottlenecks. Some games are bottlenecked by the application itself, right, doing your usual AI physics, et cetera. In this case, it was probably limited by the graphics processing just to avoid rendering possibly simple geometry. Don't do anything more complicated than simple sorting, culling, and static LOD management. The other area where it could be bottlenecked on the CPU is in the driver or sometimes in the API. In these cases, then there's something very strange going on because we're having to do a lot of work to run the graphics chip. The biggest cause of a driver bottleneck is a great number of batches per frame and we covered the details before. So keep track of those. Also sometimes you may be doing begin and end scene for every draw call. That's not good, either. So just check that you're doing something reasonable with the API. And it should be noted that most graphics applications these days are CPU limited, even with your newer CPUs.

CPU瓶颈：解决方法

CPU bottlenecks: Solutions

- 使用CPU模板（profilers）（例如Intel的VTune）
 - 驱动程序可能大部分时间都处于空闲
 - 寻找汇编程序的空循环或者使用NVPerfHUD
- 尽可能地增加批处理的尺寸
 - 所需要的代价是在GPU里！
- 对于渲染
 - 宁愿GPU负荷重一点，也要减轻CPU负荷
 - 避免复杂的（但开销更大的）CPU算法设计以降低渲染负载



#NVIDIA

The easiest way to check CPU limitations is to use V-tune, for example. V-tune should let you know how much time you're spending in the application, in the API, and/or in the driver. We do have something called NVPerfHud which should help you out, not only for this but for a couple of other things, and I'll show you an example later. A solution for the batch size problem is to try to increase the number of polygons per batch and, therefore, reduce the number of batches. We talked a little bit of how to do this using texture pages and we're going to go into some detail about how to optimize the other parts of the GPU but I want you to remember that whatever I'm saying about the optimizations you do for the GPU, don't do them if they're going to impact or increase the number of batches dramatically. So when you're optimizing for batches, do it at the expense of the GPU.

AGP传输瓶颈

AGP transfer bottlenecks

- 不像AGP4X的瓶颈
 - AGP8X在此
- 太多的数据通过AGP总线：
 - 无用的数据
 - 解决方法：减小顶点尺寸
 - 太多动态顶点
 - 解决方法：使用顶点着色器动画来代替
 - 动态数据的低效率管理
 - 解决方法：使用正确的API调用和标记
 - 视频缓存的溢出
 - 解决方法：确保帧缓冲、纹理和静态顶点缓冲大小适合视频缓存的尺寸



#VIDIA.

Now, let's step through all the stages of the pipeline after the CPU. The first thing we hit is the AGP bus. Like I said, most of the time, you're not going to be limited by this. The AGP8X bus transfers data at a rate of 2 gigabytes per second. That's plenty of bandwidth to transfer tens of millions of vertices per second. If you do find that you're limited by this, you can reduce the vertex size. We do have vertex formats that are 16 bits instead of 32 bits per component and that should help quite a bit. Also, if you find that you are using too many dynamic vertices, then chances are you're not using vertex Shaders enough. You can do complicated animations on the vertex shader such as skinning and morph targets and you should be doing this instead of manipulating geometry on the CPU. Next, if you do need to have a lot of dynamic data, you should be managing this data with the proper API calls and with the proper flags. We'll go into a little more detail of the correct flags for vertex buffers. And, finally, just make sure that everything is fitting into video memory. So what can happen is that sometimes textures can spill over into AGP memory and then, at some point, your scene hits this texture that's in AGP memory and you're slowing down because you don't have the 20 gigabytes per second to video memory access, you have the 2 gigabytes per second to AGP memory access and the GPU is then sitting idle waiting for this data to be transferred.

AGP传输瓶颈

AGP transfer bottlenecks

- 使用不合适的格式传输数据：
 - 顶点尺寸应该是32字节的整数倍
 - 解决方法：调节定点尺寸达到32字节的整数倍：
 - 压缩各个分量并使用顶点着色器进行解压缩
 - 增加数据达到更高的一个倍数
 - 非有序地使用顶点 (pre-TnL cache)
 - 解决方法：重新对顶点排序以保证顺序使用
 - 使用NVTriStrip



#NVIDIA

When transferring vertices across the AGP bus, you want to make sure that they are in 32-byte multiples. So you're going to have 32, 64, et cetera and ensure that they're being transferred sequentially. We do have a triangle stripper called NVTriStrip that you can download from our site and this will actually reorder the vertices to be in sequential use and make optimal use of the post-T&L vertex cache.

优化几何传输

Optimizing geometry transfer

- 静态几何:
 - 创建一个只写的顶点缓冲区并且只往里面写一次
- 动态几何:
 - 创建一个动态的顶点缓冲区
 - 在帧开始的时候锁定DISCARD
 - 然后附加NOOVERWRITE直到填满
 - 更频繁地使用NOOVERWRITE而不是DISCARD
 - 每一个DISCARD要么需要更多时间要么需要更多存储器空间
 - 因此NOOVERWRITE应该是最常用的
 - 决不要在未标记的情况下使用
- 半动态几何:
 - 为程序或者按需装载的几何体
 - 一次锁定，在很多帧里面使用
 - 可以尝试静态和动态两种方式



#NVIDIA

When managing your vertex buffers, you want to create them with a write-only flag if they're going to be static. If you're doing dynamic geometry, then you want to create a dynamic buffer with the dynamic flag, you want to lock with discard and overwrites and the rule of thumb for this is that an overwrite doesn't hit us as bad as a discard. Discard means that the driver has to allocate a new vertex buffer somewhere in order for you to write to it while the GPU is actually using the previous vertex buffer at that point in time. And the takeaway from both of these rules is that you should never not use a flag. Always, always, always tell us how you're going to be managing your geometry. For semi-dynamic geometry, you can try a bit of both. What we've found is that, if you only update your vertex buffer once per frame, for example, you can create it as static and do a discard on it. And that's usually pretty good because it makes it reside in video memory. If you're creating a dynamic vertex buffer, it'll typically reside in AGP memory so just bear in mind where you want to store it and how often you're going to modify it.

顶点变换瓶颈

Vertex transform bottlenecks

- 一般情况下顶点变换不太可能会成为瓶颈
 - 除非你的场景每帧有一百万个三角形 (Cool!)
 - 又或者顶点着色器 (vertex shader) 用爆了 (Cool!)
 - >例如采用了128条以上的vertex shader指令
- 太多顶点
- 解决方案: 使用LOD (level of detail, 细节等级)
 - 但是: 由于GPU有强大的顶点处理能力, 很少会因此造成性能下降。
 - 因此: 请尽量地在CPU上执行拆分场景模型的侦测、运算。
 - 两到三级的静态LOD效果就相当不错了



#VIDIA.

Next in the pipeline we have our vertex engine. Now, for the GeForce FX 5900 ultra, we can transform a whole lot of vertices -- we're talking in the order of hundreds of millions of vertices per second. So I really doubt that you're going to be bottlenecked by your Shader instructions. Usually you could be bottlenecked here because your vertices are not in strip order, meaning that you make very poor use of the vertex cache. But let's say that you have too many vertices in general, what you could do is use static LODs. Basically drop off the LODs in the distance but, don't do any smart things in the CPU to do dynamic subdivision, for example. It's possible to do but you don't want to be doing it for all your objects.

造成顶点变换瓶颈的原因

Vertex transform bottleneck causes

- 每个顶点上进行了**太多的运算**:
 - 在顶点光照方面, 进行了大量复杂光照或者光照方式运算 (指本地视点角度的情况下, **local viewer**)
 - 方向光(**Directional**) < 点光源(**point**) < 聚光灯(**spot**)
 - 启用了**Texgen**又或者纹理矩阵不一致
 - 顶点着色器程序存在:
 - 大量的指令
 - 大量的循环叠代或者分支
 - **Post-TnL vertex cache**利用率低下
 - 使用**nvTriStrip**



#NVIDIA

The next thing in the vertex engine is the number of instructions. Now, as I said before, you can transform a whole lot of vertices and, obviously, the more instructions that you execute per vertex, the slower you're going to go. Sometimes you are applying lights that, for example, will not actually light that object. So you may want to be toggling these, switching between different Shaders to ensure that you're optimally processing that vertex. The next thing is that you may be using a lot of the newer instructions such as looping or branching. Looping itself is not particularly bad. Branching starts to get worse. I mean, just like any processing unit, the more you branch, the less efficient you get and it's a lot harder for us to do branch prediction on the GPU. So try to limit it, for now, though in the future, this is going to be a faster approach to-- especially for increasing batching.

顶点变换瓶颈的解决方案

Vertex transform bottleneck solutions

- 重新排列顶点，使之能以连续的方式在**PostTnL cache**中使用。
 - **NVTriStrip**
- 把**per-object**的计算脱离出**GPU shader**
 - 在**CPU**上进行计算并保存为程序的常量
- 透过采用复杂指令以及向量操作减少指令的数量
 - 或者使用**Cg/HLSL**等高级实时**shader**语言
- 对每条**mov**指令需要特别的留意
 - 或者使用**Cg/HLSL**等高级实时**shader**语言
- 考虑**level of details**式的**shader**程序
 - 远处的对象真的也需要**4-bone skinning**吗？
- 考虑把基于**per-vertex**的工作转移到基于**per-fragment**上来
- 强制提高屏幕分辨率或者使用抗失真（**anti-aliasing**，例如**FSAA**、运动模糊、各向异性过滤等）！



#VIDIA.

Sometimes, when we write Shaders, we write them for convenience. In other words, we just put in a few constants, maybe multiply them in the Shader itself, instead of actually calculating them on the CPU. So these kinds of compilations can generally be removed; however, remember that you shouldn't do this if it's going to impact the number of batches that you're sending across. So, for example, it might be a good idea, if you have lots of batches, to concatenate them and do some per object calculations on the GPU but, if you're not suffering from batching performance and you're suffering from vertex Shader performance, then take out the instructions, basically just constant instructions, and leave them to the CPU. Next, recall that the vertex Shader is a 4D unit. In other words, it's got four units that operate on a vector simultaneously and it's got lots of complex instructions that have been heavily optimized by our architecture team. So you can use a combination of those complex instructions and also remember that you have vector operations that can be swizzled so you can have inputs that map to different source inputs and then output to maybe just one, two, three or even four components of the vector without having to start moving data around yourself. And that brings us to the next point which is we very, very rarely see a need for a move instruction so every time you see one, make sure it's for a very good reason. And, finally, there's a lot of work that's probably not necessary for objects that are really far away. So before we talked about reducing the number of vertices on objects that have lower levels of detail for when they're in the distance. Well, you can do the same thing with shaders so one example is, if you have a character that's skinned with four bones, well, does he really need four bones when he uses only a few pixels on the screen? Chances are not and you can drop that to two or maybe even one bone. And another thing that you could do is, instead of reducing the bottleneck here, you can just increase the bottleneck somewhere else and you can either increase the resolution or enable antialiasing or anisotropic filtering.

坐标设定 (Setup) 操作瓶颈 Setup bottleneck

- 几乎从来就不是瓶颈
 - 除了专门针对这类操作的性能测试程序外
- 该操作的性能会受以下情况影响：
 - 三角形的数量多寡
 - 需要进行光栅化处理的每个三角形顶点属性数量
- 想提高Setup操作性能的话，请：
 - 减少合并为真正三角形的比率
 - **But only if that ratio is substantial (> 1 to 5)**



#NVIDIA

Next module in the pipeline is the setup unit and this is where the vertices get made into triangles or quads. This is almost never a bottleneck so you probably don't have to worry about it. It's influenced by the number of triangles or the attributes that each triangle has. In other words, it's got to set up each of the colors, the texture coordinates or whatever other data you're passing around so the more of that data you have, the longer it'll take to set up that triangle but, typically, this is insubstantial.

光栅化瓶颈

Rasterization bottlenecks

- 有大量**z-cull**三角形时才可能出现瓶颈 **triangles**
 - 罕见
- 受以下因素影响其速度：
 - 三角形的数量
 - 三角形的大小尺寸



#VIDIA.

Next is the rasterization module which breaks up the triangles into pixels and, again, this is not a module that typically causes a bottleneck. The only time that this is going to be a bottleneck is if you have a lot of non-visible triangles that are being Z culled. Basically while this triangle's being Z culled, it's introducing lots of bubbles in the pipeline and effectively generating no-ops for all the stages downstream from this rasterization module.

GPU瓶颈 – 片段着色器

GPU bottlenecks – fragment shader

- 在**过去的**架构体系中，由于着色器程序的功能固定、搭配简单，其执行性能和流水线的峰值性能相当。
- 在**NV1X (DirectX 7)**中，使用更多的**通用组合器(general combiner)**会造成片段着色器的性能下降，但是一般来说不会构成瓶颈
- 在**NV2X (DirectX 8)**中，片段着色方面引入了更复杂、需要消耗大量运算资源的**片段着色器(fragment shader)**模式。
- **NV3X (CineFX / DirectX 9)**能够跑**512条指令**（**OpenGL**中是**2048条**）的片段着色器程序。
 - 较长的片段着色器程序会造成瓶颈



#NVIDIA

So now we get into the really interesting part, the fragment shader, and this module is much more complicated than it was in the past. What this means

is that there is a lot more opportunity to make it **go** slow. Now, with the CINE-FX architecture, basically all of the GeForce FX family, you can run 512 instructions under DX-9 or 2048 in Open GL. And obviously the longer the shader, the slower it goes and it's very easy to write a long fragment shader, especially when using high level languages. Again, when we're writing shaders, we don't necessarily think about the performance of them until we start running them and you are going to have to write them liberally at first and then go back and optimize them and see where you can start squeezing more performance out of them. I'll go through some points that will cover these topics.

GPU瓶颈 – 片段着色器：原因与解决方案

GPU bottlenecks – fragment shader: Causes and solutions

- 太多的片段（**fragments**，或者说待处理像素）
- 解决方案：
 - 尽量以**front-to-back**的顺序来渲染场景
 - 考虑在第一个**Pass**的时候只渲染**Z-buffer**(关闭所有特效)
 - 让你在后续的**pass**中只渲染可见的**fragment**
 - 但是: 你依然需要牺牲顶点性能来提高**fragment**性能
 - 因此: 请不要对简单的**fragment shader**进行这样的操作
 - 提醒一下，这样的方式能节省帧缓存带宽



#VIDIA.

In the first place, you might have too many pixels and this will cause a bottleneck in the fragment shader. What you can do here is to try to use our Z cull engine more efficiently by rendering roughly front to back. Another thing that you can do if you have complex shaders or lots of multi-pass shading is that you want to do a Z only first pass. When you're doing a complex, single pass shader, you want to render Z only, effectively mask out your colors and you can do this using a render state. This will actually let us run incredibly fast because we know that you're not writing any color, so we can really kick the engine into turbo and spit out these pixels and do really fast Z checks. Now, what happens is, that next time when you actually render the real pixels, you will have already written the entire depth surface and you will only be lighting the pixels that are visible. So this is a way of decoupling geometric or your scene complexity from your pixel complexity. Now, what happens if you have multi-pass is you can do the same thing. Instead of writing Z only first, you just write the first pass, including Z and then, for subsequent passes, make sure you turn off Z writes but keep Z checks on.

GPU瓶颈 – 片段着色器：原因与解决方案

GPU bottlenecks – fragment shader: Causes and solutions

- 每个fragment进行了大量的运算
- 解决方案:
 - 利用复杂指令、向量操作、并行发射 (co-issuing, RGB+Alpha) 来减少指令数量
 - 混合使用texture和combiner指令 (这些指令能并行运行)
 - 使用配对的combiner指令
 - 使用配对的 (简单) texture指令
 - 使用alpha混合器来帮助提高性能:
 - SRCCOLOR*SRCALPHA for modulating in the dot3 result
 - SRCCOLOR*SRCCOLOR for a free squaring
 - 考虑使用LOD shader
 - 对于远处的对象关闭detail map计算
 - 考虑把per-fragment的计算迁移到 per-vertex



#VIDIA.

Another cause of this stage being a bottleneck is that you are just running too many instructions per pixel. Just like the vertex shader, this unit operates on 4-D vectors. You can make it actually run on a 3-D vector and a scaler, or what we consider Alpha. This is called co-issuing. You can actually save quite a number of cycles by doing co-issuing and it basically happens automatically as long as you're not using the alpha, the result of the scaler operation in your vector operation. Next is that you can optimize by doing even texture lookups and even map operations. On our hardware, you get two types of lookups and two map operations in a single cycle. If you're not running all these modules, if you're not using all these modules at the same time, then you're basically just going to be wasting their cycles. Next thing is that sometimes you might be able to get away with doing multiplication in the alpha blend unit so there's a final multiply out the back door. Do note that this happens at 32 bits per pixel. In other words, it doesn't happen with floating point formats. And, again, you can do shader level of detail, like we discussed for the vertex engine. Can you really see that the character that's using up five pixels of the screen is truly bump mapped or can you even see that there are eight lights on him? Chances are you can't so just, kill certain instructions, load a different Shader when the object is far away. Another thing that we talked about in the vertex engine was to move constant to work. In other words, by sending them to the CPU. Now, this also applies to the pixel Shader because, basically, if something is going to be linear, it can be done in the vertex engine. If something's going to be constant, it can be put into a constant register. So, consider moving stuff from the actual pixel shader to the other units.

GeForce FX的fragment shader优化 I

GeForce FX fragment shader optimisations

- 额外的性能最佳化指引：：
 - 尽可能地使用**fx12精度**指令。
 - 之前所有版本的DirectX都不过是**9-bit精度**
 - 尽可能地使用**fp16**指令
 - 运行效果要比传统的色彩混合好得多
 - 使用**_pp**指令**modifier**
 - 在**HLSL**中打开**ps_2_a**描述(**profile**)开关
 - 尽可能地减少临时存取（寄存器数量）
 - 在能接受的时候使用**16-bit寄存器**(大多数情况下都可如此)
 - 重复使用寄存器并且使用各寄存器的所有分量(**component**) (**swizzling**方式的操作是**free**的)



#VIDIA.

As I said before, the CINE-FX architecture is more complicated than previous ones and we do have a variety of precisions within the fragment shader. The lowest precision is fixed point 12 base and it basically gives you a range of minus two to two. So this is pretty good, considering that, for all previous Direct X versions, the maximum precision was about 9 bits so 12 bit fixed should give you enough to do most of your usual lighting calculations. The next step up from there is floating point 16. It's actually what we call the half precision format and this is great for slightly greater range. It works fine for anything that was previously done with 9 bit precision, obviously, and it's accessible under the DX-9 assembly instructions when using the underscore PP modifier. In HLSL, you can access the half data type using PS2_A and this is also optimized for the CINE-FX architecture. Now, because HLSL or CG makes it very easy to write long, complicated pixel shaders, it's also very easy to dramatically increase the number of temporary registers that you'll be using. So the point is to use half precision registers if possible, i.e., the half data type or the FP16 float registers. And also try to reuse as many registers as possible so don't load up 16 textures into 16 different registers and then use each one individually when you could have actually just loaded two textures at a time into two registers, done a math operation, saved it into one of those whose value was discarded before and then reuse another two registers for the next set of texture loads. The other thing that you can do is that we do have arbitrary swizzling so you can, for example, store two 2-D vectors in a single 4-D vector because this is a 4-D engine and just swizzling, masking two components, is absolutely free.

GeForce FX的fragment shader优化 II

GeForce FX fragment shader optimisations

- 根据不同的操作选择合适的单元: CPU, Vertex, Pixel
 - CPU for constants, Vertex for linear calculations
- 采用最低版本的pixel shader
- 只有以下情况你才需要每个周期那么多的pixel shader
 - 把他们用于有趣的视觉特效
 - 例如Per-pixel 凹凸映射(bump map)和反射(reflections)
- 用更多的时间去努力减少不必要的精度和计算



#VIDIA.

In the end, what we're talking about is to choose the appropriate unit for the type of data that you're calculating. Constant data can definitely go into the CPU and then be transferred to the vertex and pixel engines as constant registers. Vertex engines used for data that's going to be interpolated across a triangle and the pixel engine is going to be used for data that's going to vary at every pixel. Do try to use the lowest pixel shader version available to you. Basically, what this means is, if you're only doing 4 texture lookups and 8 map operations and you're not doing a dependent texture read, then use PS-1.1 up to 1.3 version. If you're doing some dependent operations and you can fit within 1.4 shader, then use that and, if you're doing lots of dependent operations, then use, only then use pixel shader 2.0. This also guarantees that, for previous pixel shader versions, 1.X, that you're using the low precision register types so that actually saves storage and cycles. The same as you get a certain number of batches per frame, you can also think of getting the same number of pixel Shader cycles per frame. So you've got to basically choose where you're going to spend your cycles. Do effects that are going to distinguish your game visually, stylistically. Don't do effects that have a minor difference on a pixel and they're consuming several cycles to do so.

GPU瓶颈 – 纹理：原因与解决方案

GPU bottlenecks – texture: Causes and solutions

- 纹理太“肥”了：
 - 纹理cache溢出：出现大量的cache命中失败
 - 视频内存溢出：必须透过AGP来拾取纹理
- 解决办法：
 - 纹理分辨率应该适得其所即可，尽量不要用更高分辨率的纹理
 - 避免“昂贵”的内部格式
 - CineFX可以使用浮点4xfp16和4xfp32格式
 - 压缩 纹理：
 - 把单色通道收归到alpha里
 - 对于像环境贴图和阴影贴图的纹理，尽可能使用16-bit色彩
 - 使用DXT压缩，提醒一下，NVIDIA现在的GPU在DXT1方面的品质相当出色



#NVIDIA

Okay. As part of the input to the fragment engine is the texture unit. The texture unit basically consists of the cache and a decoder for compressed textures. One of the problems that you can run into is that you are thrashing the texture cache a lot or you're loading textures from AGP memory. I mentioned how you can spill over from video memory so one of the things to do for the latter is to use compressed textures. Use 16-bit colors whenever you can. Use DX T-1 compression over the other compression formats but, if that's not good enough quality, then go on to DX T-3 and T-5. Now, textures should also be using MIP maps and when choosing the highest MIP level, basically choose it depending on how much of that texture is going to be on screen at a certain screen resolution. So you don't really want to be drawing a character with 2048 x 2048 textures when, at most they're using a 512 x 128 area of the screen. Now, with the new capabilities that CINE FX offers is that you can choose to store floating point data. Now, floating point data is expensive for two reasons. First, it uses a lot of bandwidth to video memory, whether it be through textures or through render targets. Secondly, it costs a certain amount to decode these floating point surfaces. So use them sparingly. Use them only when needed and you only really need them for high dynamic range and actually the next slide here shows you the format range and precision that they run at.

纹理格式, 范围, 精度

Texture format, range, and precision

Format	Range	Precision	Blend	Render	Texture	Filter
A8R8G8B8	2.4 dB	8 bit	Y	Y	Y	Y
RGBE8	76.8 dB	8 bit	N	Y	Y	N
D3DFMT_R16G16	4.9 dB	16 bit	N	N	Y	Y
RGBA16F	9.3 dB	11 bit	N	Y	Y	N
RGBA32F	76.8 dB	24 bit	N	Y	Y	N



#NVIDIA

So your typical 32-bit ARGB format gives you 256 levels, it's not much but it supports all the blend, render, texture and filter operations. ARGB-8 is basically the same formats but you're using the alpha as an 8-bit exponential that's shared across the three color components. Now, this gives you a lot more range; however, even though it's a standard format, blending and filtering doesn't really work on that because the exponent can't just be added. The next format is 16-bit high low format. Basically 16 bits integer red, 16-bit integer green making up a 32-bit slot. This gives you slightly more range but the problem is that you can't render to it. Because you can't render to it, you can't blend to it. However, they're great for textures and, in effect, they're actually pretty good for storing really high precision vectors. For example, you can use two of these texture surfaces to make a 4-D 16-bit integer surface that is filtered and running optimally on a hardware. Now, we'll get to the floating point formats. First one is the half-precision. This has greater range than the previous integer formats but you can't blend to it and you can't filter it. And then the 32-bit version of that has the greatest range of all but, again, you can't blend to it and you can't filter to it. In the future, this is likely to change as it becomes much more widespread and more useful.

GPU瓶颈 – 纹理贴图：原因与解决方案

GPU bottlenecks – texture: Causes and solutions

- 纹理**cache**利用率低下：大量的**cache**命中失败
- 解决方案：
 - 局部式的纹理存取
 - 相关式纹理查表（**dependent texture look-up**）
 - 使用 **mipmapping**:
 - 避免负值 **LOD偏移值** 来做纹理锐化：纹理**cache**是为标准**LOD**设计的
 - 锐化通常会导致远方出现失真(**aliasing**)
 - 最好使用各向异性过滤(**anisotropic filtering**)做锐化
 - 小心**边长为非二次方的纹理**（**non-power of 2 textures**）
 - **Cache**利用率通常都不如等边的纹理



#VIDIA.

What happens if the texture cache is under-utilized? Well, chances are that you're not using the texture cache well. The easiest way to break the texture cache is to do lots of dependent texture reads and what I mean by dependent texture reads is to generate a texture coordinate in the fragment shader and then go fetch it instead of the interpolated texture coordinate. Let's say you're actually generating texture coordinate from a previous texture read, so, in this case, we've got some sort of random data coming in, doing some math operations and generating brand new texture coordinates. This is going to start fetching texels from areas that are very different in your actual texture map, which is going to cause the cache to be thrashed often. A good way to start avoiding this is to MIP mapping because, basically, calculating the gradient of your dependent texture reads enables us to fetch the texels from a smaller MIP map. But do avoid using negative LOD bias. Negative LOD bias basically sharpens the image but what happens is that it looks sharper for a static image but it actually causes aliasing on the moving image. So, on a screen shot, it might look nicer but, when you actually play the game, you'll see sparklies going around. Instead of doing negative LOD bias, if you do want to sharpen, you should use anisotropic filtering only for those textures that really need it. So what anisotropic filtering does is to sample more in a given direction. We choose the direction by calculating how much texel area is used by a pixel. So if you're using more texels on the X axis than on the Y axis, then we're going to do more filtering on the X axis. And we can have up to an 8:1 ratio of filtering. We can take eight times as many samples in the X direction than in the Y direction. And lastly you should beware of non-power-of-2 textures. They tend to have worse caching behavior than power-of-2 sizes, especially when doing dependent texture reads.

GPU瓶颈 – 纹理贴图：原因与解决方案

GPU bottlenecks – texture: Causes and solutions

- 每次查表进行了太多的取样
 - 三线性（**Trilinear filtering**）会让填充率减半
 - 各向异性过滤更甚之
 - 取决于各向异性的级别
 - 硬件能够“精明”地处理各向异性过滤，但是你必须购买具备这样功能的硬件
- 解决方案：
 - 只在适当的地方使用**trilinear**或者**anisotropic**过滤：
 - 一般来说，只有**diffuse map**真正需要
 - **Light map**的分辨率太低难以察觉**mipmap**分界线而需要不大
 - **Environment map**的应用场合经常出现扭曲，也没多大的必要
 - 降低 **anisotropy**的比率
 - 一般来说使用了**anisotropic**过滤的话可以减少 **trilinear**过滤

Another texture unit problem could be filtering so tri-linear filtering literally cuts your fill rate in half. Good way to avoid that penalty is to only enable tri-linear filtering on textures that really need it. Light maps, for example, are typically very low resolution and you won't be able to see the MIP line or you might not even hit a lower MIP level until very, very far away, at which point it hardly matters. So there are certain textures that you don't need tri-linear filtering for and, equivalently, there are some textures that you don't need anisotropic filtering for. Pick and choose what type of filtering you do, minimize expensive ones as much as possible.

快速的纹理上载

Fast texture uploads

- 使用API规范的资源管理方式，而不要使用你自己的方式
 - 请相信大多数贴图操作时候run-time和驱动程序的表现
- 对于真正的动态纹理：
 - 使用D3DUSAGE_DYNAMIC和 D3DPOOL_DEFAULT 来创建
 - 用D3DLOCK_DISCARD锁定
 - 永远都不要读取这样的纹理！



#NVIDIA

Now, what happens if you're running out of texture space or for some reason you're updating a texture dynamically? Well, what you want to be doing is, again, using the flags that we talked about for vertex buffers. You want to be creating a texture with a dynamic flag whenever you're modifying it and you're going to be locking them with discards but you never, ever, ever want to read a texture. So lock to update it but don't read it. Just discard it. The driver will start up a new area in memory in case the GPU was busy with the previous texture, write to that and then discard the previous memory area once the GPU is done using it.

GPU瓶颈 – 帧缓存：原因与解决方案

GPU bottlenecks – frame buffer: Causes and solutions

- 太多针对frame buffer的读写
- 解决方案:
 - 关闭Z写 (Z writes) :
 - 对于多Pass方式的渲染，由于能够在第一个pass的时候就完成z-buffer的渲染，因此后续的pass就可以关闭掉无用的Z write。
 - 对于alpha-混合的几何对象(例如particles)也可以关闭z-write
 - 但是, 不要只对个别的色彩通道进行mask off操作:
 - 这实际上会导致性能下降，因为GPU必须在写有掩码的色彩通道时，必须先从帧缓存中读取这些通道来
 - 使用alpha test (除非你把所有的色彩通道做了mask)
 - 对浮点帧缓存的使用提出质疑
 - 这需要非常非常多的带宽



#VIDIA.

Now we're getting into the final stages of the graphics pipeline and here we're basically talking about reads and writes to the frame buffer. What could be happening is that your read/modify/writing to the Z buffer, which is not always necessary. For example, for multi-pass rendering, after the first pass, you can turn off Z write and instead only do compares. Also, for alpha-blended geometry, you may not need to do a Z write. Now, DX-9 allows you to mask off color channels but if you only mask off some color channels, what happens is that the GPU has to go read the frame buffer, read the masked channels, combine them with the non-masked channels and then write that pixel back again so you either want to turn color writes on or off but don't do partial writes. And, lastly, again on the subject of floating point buffers, you really want to question their use because they use up a tremendous amount of memory bandwidth.

GPU瓶颈 – 帧缓存：原因与解决方案

GPU bottlenecks – frame buffer: Causes and solutions

● 解决方案 (继续):

- 如果没有用到**stencil**的话尽量使用 **16-bit Z 深度**
 - 许多室内场景使用**16bit Z-buffer**的效果就不错了
- 减少**render-to-texture target**的使用和尺寸
 - **Cube map**和**shadow map**都能够在使用低分辨率和**16 bit**色彩深度的情况下提供不错的视觉效果。
 - 在处理反射的时候尝试采用**hemisphere map**来替代**cube-map**
 - 文件尺寸能够比效果接近的**cube map**更小
 - 更少的**render target**切换
 - 多次重用**render target texture**来减少内存覆盖区
 - **Separate Z buffers** or else **Z-cull will be invalidated**



#NVIDIA

Other tricks that you may be able to get away with is using 16-bit Z buffers. Typically in an interior scene, you don't have a large Z range so this might be good enough for you. Don't trade this with W buffers because our hardware is fully optimized for Z buffers, not for W buffers. Otherwise you're going to lose things like Z cull, for example, and that you definitely don't want to lose. For techniques where you're rendering to a texture, like rendering the environment map for some object, then you want to reduce the resolution of these environment maps or you could even reduce the color depth of these environment maps. Let's say you're mapping this on a sphere, especially on a slightly diffuse or very bumpy sphere, then you won't be able to see the texture detail and the resolution will be useless at that point. So really consider what quality you're looking for, depending on what it's being used for. Now, we talked about spilling over the video memory to AGP memory. Well, sometimes what can happen is that, if you're rendering to off-screen surfaces, you may be consuming so much memory that the driver has to remove textures from video memory. In this case, you want to be re-using those render targets as much as possible; however, keep separate Z buffers because, if you are in the middle of a frame and you render to some off-screen surface without switching the Z surface as well, when you go back to draw the main render target the Z cull buffer will be invalidated and you'll lose all Z cull optimizations for the remainder of the scene.

GPU瓶颈 – 帧缓存：原因与解决方案

GPU bottlenecks – frame buffer: Causes and solutions

● 解决方案 (继续):

● 使用硬件快速路径:

● 清缓存

● Z buffer和 stencil buffer都是同一个 buffer, 因此:

- 如果你使用stencil buffer的话, 意味着Z和 stencil buffer是一并清除的
- 如果你没有使用stencil buffer, 请建立Z-only的深度表面 (例如: D24X8), 否则会影响Z清除的优化

● 在 Z-bias、 alpha test关闭以及 stencil buffer没有使用的时候, Z-cull才是最优化的。

● 尝试使用DirectX 9新的常数色彩混合来代替全屏的quad做染色特效(tinting effect)

● D3DRS_BLENDFACTOR

● OpenGL 1.2以上版本均支持该操作



#VIDIA.

Lastly, a few general tips to keep in mind. This has been true for all hardware for awhile now. Buffer clears are very fast but do try to keep buffer clears together because the Z surface actually includes the stencil surface, then clearing them together avoids spending time clearing them separately and that's not going to help. We can actually clear Z very, very fast and, if we are able to clear a stencil as well at the same time, then it's practically for free. The Z cull buffer, as well, is particularly fast when you're not using Z bias or alpha testing or the stencil buffer. So try to do as much geometry without using those three things and, lastly, you can actually do color blends for free if, for example, you're rendering the last pass of a multi-pass because you're already doing blending.

总结

- 现在**GPU**可**编程流水线**, 这意味着更多潜在的瓶颈和**更复杂的优化**
- 优化的目标是在渲染场景中每一部分时都能**让流水线的每个阶段保持不空闲(包括CPU)**
- 明白场景中什么因素制约了你的程序性能
 - **Skybox**(天空)有可能是纹理方面的瓶颈
 - **Skinned**、**dot3**角色有可能是受到传输或者变换方面的瓶颈制约
- **剖析低效率的原因, 免费地增加“特效”**
 - 使用了“昂贵”**fragment shader**特效的对象通常能免费地或者以很低的代价执行看来“昂贵”的**vertex shader**特效。



#VIDIA.

We've seen that current generation hardware is quite complicated. We have stepped into the realm of highly programmable architectures which means that it's getting much easier to do bad things and it's also getting harder to identify what you're doing incorrectly. So, hopefully, this presentation has given you a good place to start when optimizing your graphics engine and remember that, because rendering a frame is a pipeline from the CPU to the last stage of the GPU, you want to keep each stage as busy as the slowest stage. Conversely, you want to reduce the load of the slowest stage as much as possible and you have to understand that each part of the scene is going to have a different bottleneck. If you're doing multi-pass rendering then each pass of that multi-pass rendering might have a different bottleneck. And, lastly, just try to get more quality out of the scene if you're not able to reduce the workload.

Questions, comments, feedback?

- Koji Ashida <kashida@nvidia.com>

- <http://developer.nvidia.com>



#VIDIA.

Thank you very much. I hope this was useful and, in particular, that you will be able to make a significant difference to your engine. You can download NVPerfHUD from our developer site.