# Performance Tools

**Raul Aguaviva and Jeff Kiel (NVIDIA Corporation)**

# Performance Tools Agenda

- **GPU architecture at a glance**
- **Intel VTune: Code Profiling**
- **NVGLExpert: OpenGL API Assistance**
- **NVShaderPerf: Shader Performance**
- **NVPerfKit: Driver and GPU Performance Data**
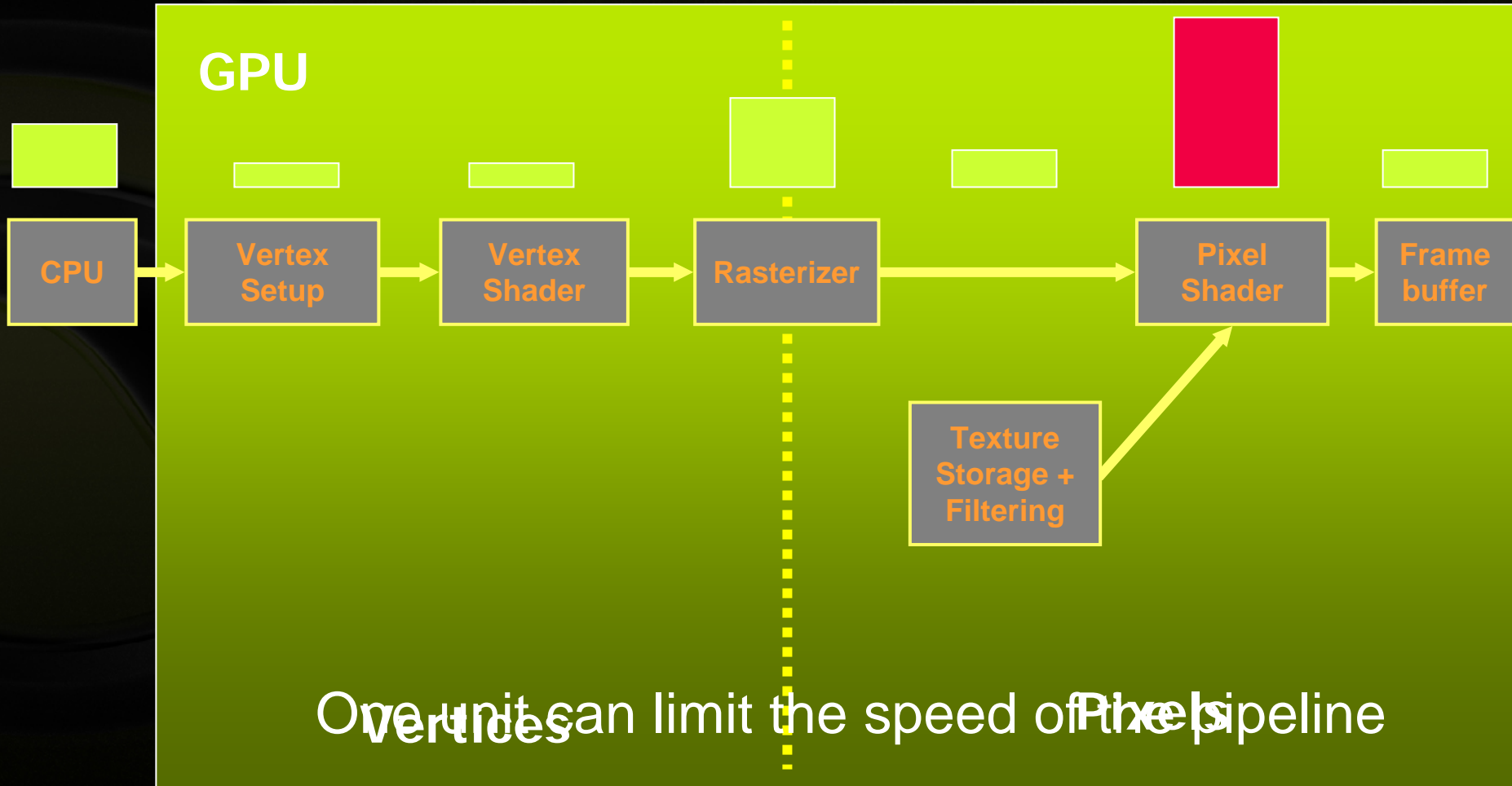- **NVPerfHUD: Interactive Performance Analysis**
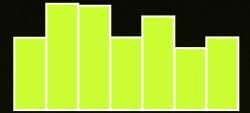
# GPU architecture at a glance

- **Pipelined architecture**
  - Each unit needs the data from the previous unit to do its job
- **Bottleneck identification and elimination**
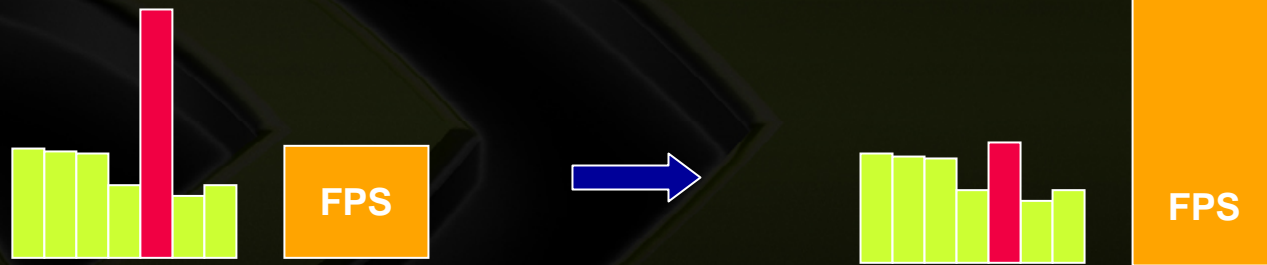- **Balancing the pipeline**

# GPU Pipelined Architecture (simplified view)

**GPU**

CPU → Vertex Setup → Vertex Shader → Rasterizer → Pixel Shader → Frame buffer

Texture Storage + Filtering → Pixel Shader

One unit can limit the speed of the pipeline

vertices                    Pixels
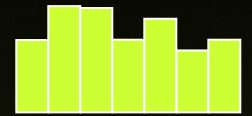
# Bottleneck Identification

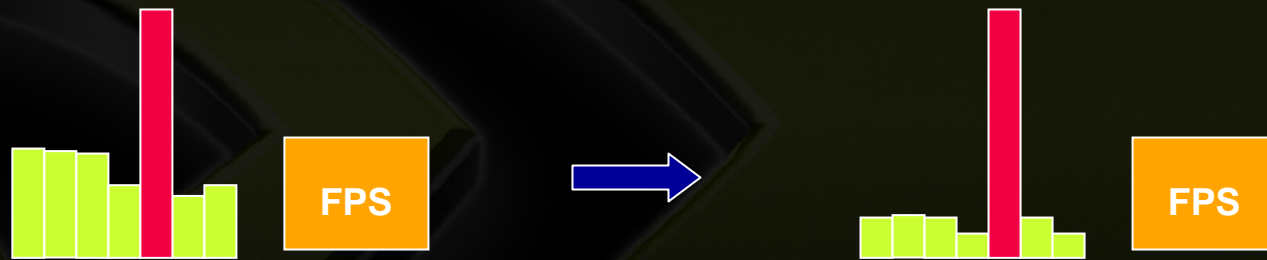- **Modify the stage itself**
  - **By decreasing its workload**



  - If performance/FPS improves greatly, then you know this is the bottleneck
  - Careful not to change the workload of other stages!

# Bottleneck Identification

- **Rule out the other stages**
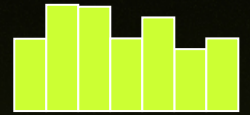  - **By giving all of them little or no work**

FPS → FPS

- **If performance doesn't change significantly, then you know this is the bottleneck**
- **Careful not to change the workload of this stage!**

# Bottleneck Identification

- **Sample counters at different points along the pipeline**
  - Use NVPerfKit and NVPerfHUD
  - How much work performed by each unit, compare to the maximum work possible

# NVGLExpert

- What is it and what does it do?
- Project status?

# What is it and what does it do?

- Helps eliminate performance issues on the CPU
- Instrumented OpenGL driver
  - Outputs information to file, console or debugger
  - Different groups and levels of information detail
- Controlled by small GUI tool
  - Windows tool sets appropriate registry entries
  - Linux tool sets environment variables
- What it can do (today)
  - Prints GL errors when the are raised
  - Indicates if the driver runs through a software fallback
  - Shows unexpected shader compile errors
  - Shows where your VBOs reside
  - Print reasons for GL_FRAMEBUFFER_UNSUPPORTED_EXT
- Feature list will grow with future drivers

# Project Status

- Will be delivered with next major driver release
- Windows and Linux
- Currently supports NV3x and NV4x architectures
- What types of things are interesting?

[NVGLExpert@nvidia.com](mailto:NVGLExpert@nvidia.com)

# NVShaderPerf

- What is NVShaderPerf?
- What's new with version 1.8?
- What's coming with version 2.0?

# NVShaderPerf

```
v2f BumpReflectVS(a2v IN,
                  uniform float4x4 WorldViewProj,
                  uniform float4x4 World,
                  uniform float4x4 ViewIT)
{
    v2f OUT;
    // Position in screen space
    OUT.Position = mul(IN.Position, WorldViewProj);
    // pass texture coordinates for fetching the normal map
    OUT.TexCoord.xyz = IN.TexCoord;
    OUT.TexCoord.w = 1.0;
    // compute the 4x4 tranform from tangent space to object space
    float3x3 TangentToObjSpace;
    // first rows are the tangent and binormal scaled by the bump scale
    TangentToObjSpace[0] = float3(IN.Tangent.x, IN.Binormal.x, IN.Normal.x);
    TangentToObjSpace[1] = float3(IN.Tangent.y, IN.Binormal.y, IN.Normal.y);
    TangentToObjSpace[2] = float3(IN.Tangent.z, IN.Binormal.z, IN.Normal.z);
    OUT.TexCoord1.x = dot(World[0].xyz, TangentToObjSpace[0]);
    OUT.TexCoord1.y = dot(World[1].xyz, TangentToObjSpace[0]);
    OUT.TexCoord1.z = dot(World[2].xyz, TangentToObjSpace[0]);
    OUT.TexCoord2.x = dot(World[0].xyz, TangentToObjSpace[1]);
    OUT.TexCoord2.y = dot(World[1].xyz, TangentToObjSpace[1]);
    OUT.TexCoord2.z = dot(World[2].xyz, TangentToObjSpace[1]);
    OUT.TexCoord3.x = dot(World[0].xyz, TangentToObjSpace[2]);
    OUT.TexCoord3.y = dot(World[1].xyz, TangentToObjSpace[2]);
    OUT.TexCoord3.z = dot(World[2].xyz, TangentToObjSpace[2]);
    float4 worldPos = mul(IN.Position, World);
    // compute the eye vector (going from shaded point to eye) in cube space
    float4 eyeVector = worldPos - ViewIT[3]; // view inv. transpose contains eye position in world space
    OUT.TexCoord1.w = eyeVector.x;
    OUT.TexCoord2.w = eyeVector.y;
    OUT.TexCoord3.w = eyeVector.z;
    return OUT;
}

///////////////////// pixel shader //////////////////

float4 BumpReflectPS(v2f IN,
                     uniform sampler2D NormalMap,
                     uniform samplerCUBE EnvironmentMap,
                     uniform float BumpScale) : COLOR
{
    // fetch the bump normal from the normal map
    float3 normal = tex2D(NormalMap, IN.TexCoord.xy).xyz * 2.0 - 1.0;
    normal = normalize(float3(normal.x * BumpScale, normal.y * BumpScale, normal.z));
    // transform the bump normal into cube space
    // then use the transformed normal and eye vector to compute a reflection vector
    // used to fetch the cube map
    // (we multiply by 2 only to increase
    float3 eyevec = float3(IN.TexCoord1.w, ...);
    float3 worldNorm;
    worldNorm.x = dot(IN.TexCoord1.xyz, normal);
    worldNorm.y = dot(IN.TexCoord2.xyz, normal);
    worldNorm.z = dot(IN.TexCoord3.xyz, normal);
    float3 lookup = reflect(eyevec, worldNorm);
    return texCUBE(EnvironmentMap, lookup);
}
```

## Inputs:
- HLSL
- GLSL (fragments)
- !!FP1.0
- !!ARBfp1.0
- PS1.x,PS2.x,PS3.x
- VS1.x,VS2.x, VS3.x
- Cg

## NVShaderPerf

## GPU Arch:
- GeForce 7800 GTX
- GeForce 6X00, FX series
- Quadro FX series

```
C:\WINDOWS\system32\cmd.exe

    dp3  r0.x, r1, r1
    rsq  r0.w, r0.x
    nrm  r0.xyz, t1
    mad  r1.xyz, r1, r0.w, r0
    nrm  r2.xyz, r1
    nrm  r1.xyz, t2
    dp3  r2.x, r2, r1
    max  r1.w, r2.x, c9.x
    pow  r0.w, r1.w, c5.x
    add  r1.w, r0.w, -c7.x
    mov  r2.w, c6.x
    add  r2.w, r2.w, -c7.x
    rcp  r2.w, r2.w
    mul_sat r2.w, r1.w, r2.w
    mad  r1.w, r2.w, c9.y, c9.z
    mul  r2.w, r2.w, r2.w
    mul  r1.w, r1.w, r2.w
    mov  r2.x, c9.w
    add  r2.w, r2.x, -c8.x
    mad  r1.w, r1.w, r2.w, c8.x
    dp3  r0.x, r0, r1
    mul  r1.w, r0.w, r1.w
    mul  r1.xyz, r0.w, c4
    add  r0.w, r0.x, c9.x
    mul  r0.w, r0.w, c10.x
    mov  r0.xyz, c0
    add  r0.xyz, r0, -c1
    ...
    mov  ...

Target: GeForce 6800 Ultra (NV40) :: Unified Compiler: v61.
Cycles: ... R Regs Max Index (0 based
Pixel throughput (assuming 1 cycle texture lookup) 304.76 M

Shader ...
Cycles: 14.00 :: R Regs Used: 2 :: R Regs Max Index (0 based
Pixel throughput (assuming 1 cycle texture lookup) 457.14 M

Shader performance using all FP32
Cycles: 21.00 :: R Regs Used: 3 :: R Regs Max Index (0 based
Pixel throughput (assuming 1 cycle texture lookup) 304.76 MB

C:\Temp\NVShaderPerf_61_77>
```

## Outputs:
- Resulting assembly code
- # of cycles
- # of temporary registers
- Pixel throughput
- Test all fp16 and all fp32
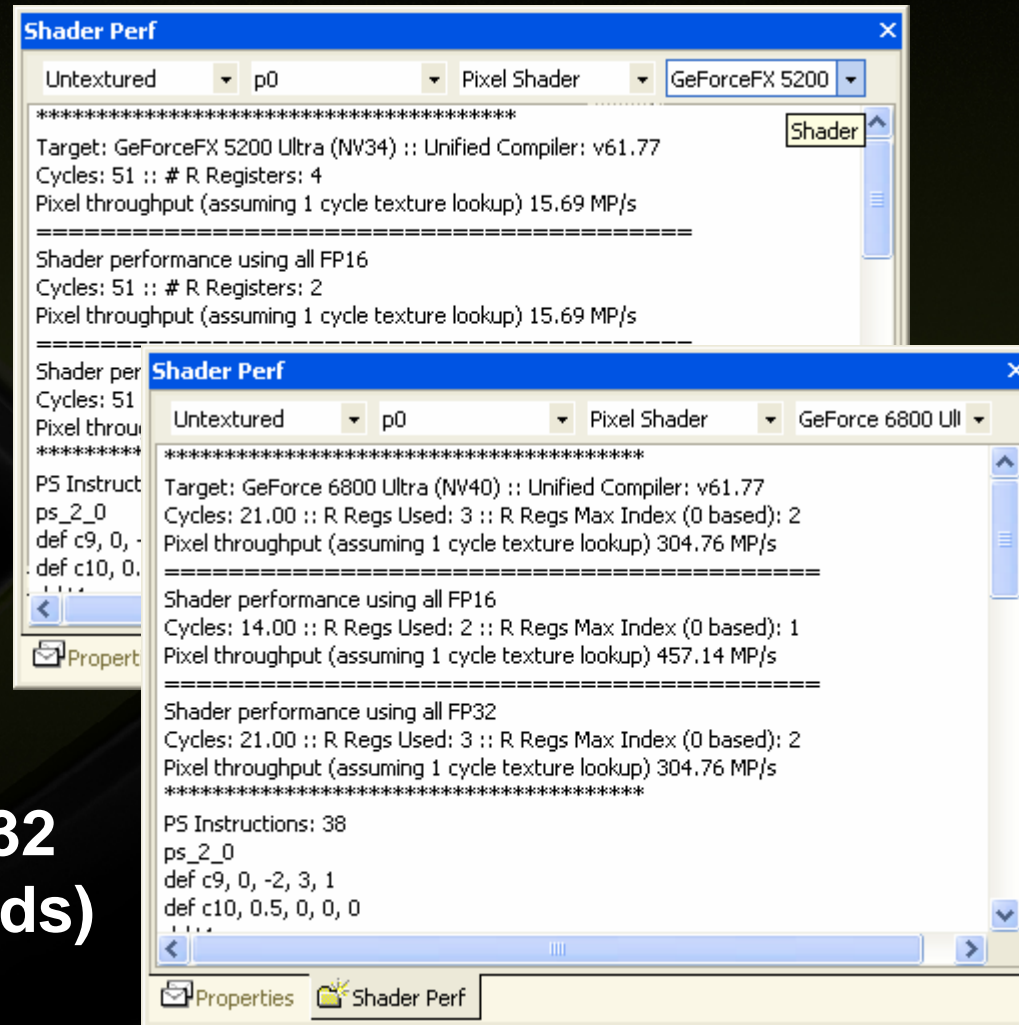
# NVShaderPerf: In your pipeline

- **Test current performance**
  - **against shader cycle budgets**
  - **test optimization opportunities**
- **Automated regression analysis**
- **Integrated in FX Composer 1.7**

# FX Composer 1.7 – Shader Perf

- **Disassembly**

- **Target GPU**

- **Driver version match**

- **Number of Cycles**

- **Number of Registers**

- **Pixel Throughput**

- **Forces all fp16 and all fp32 (gives performance bounds)**

**Shader Perf**

| Untextured | p0 | Pixel Shader | GeForceFX 5200 |

```
********************************************
Target: GeForceFX 5200 Ultra (NV34) :: Unified Compiler: v61.77
Cycles: 51 :: # R Registers: 4
Pixel throughput (assuming 1 cycle texture lookup) 15.69 MP/s
================================================
Shader performance using all FP16
Cycles: 51 :: # R Registers: 2
Pixel throughput (assuming 1 cycle texture lookup) 15.69 MP/s
================================================
Shader per
Cycles: 51
Pixel throu
********************************************
PS Instruct
ps_2_0
def c9, 0, -
def c10, 0.
```

**Shader Perf**

| Untextured | p0 | Pixel Shader | GeForce 6800 Ul |

```
********************************************
Target: GeForce 6800 Ultra (NV40) :: Unified Compiler: v61.77
Cycles: 21.00 :: R Regs Used: 3 :: R Regs Max Index (0 based): 2
Pixel throughput (assuming 1 cycle texture lookup) 304.76 MP/s
================================================
Shader performance using all FP16
Cycles: 14.00 :: R Regs Used: 2 :: R Regs Max Index (0 based): 1
Pixel throughput (assuming 1 cycle texture lookup) 457.14 MP/s
================================================
Shader performance using all FP32
Cycles: 21.00 :: R Regs Used: 3 :: R Regs Max Index (0 based): 2
Pixel throughput (assuming 1 cycle texture lookup) 304.76 MP/s
********************************************
PS Instructions: 38
ps_2_0
def c9, 0, -2, 3, 1
def c10, 0.5, 0, 0, 0
```

Properties     Shader Perf

# NVShaderPerf 1.8

- **Support for GeForce 7800 GTX and Quadro FX 4500**
- **Unified Compiler from ForceWare 77.72 driver**
- **Better support for branching performance**
  - Default computes maximum path through shader
  - Use –minbranch to compute minimum path

# NVShaderPerf 1.8

```
//////////////////////////////////////////////////////////////////////
// determine where the iris is and update normals, and lighting parameters to simulate iris geometry
//////////////////////////////////////////////////////////////////////

float3 objCoord = objFlatCoord;
float3 objBumpNormal = normalize( f3tex2D( g_eyeNermel, v2f.UVtex0 ) * 2.0 - float3( 1, 1, 1 ) );
objBumpNormal = 0.350000 * objBumpNormal + ( 1 - 0.350000 ) * objFlatNorma
half3 diffuseCol = h3tex2D( g_irisWhiteMap, v2f.UVtex0 );
float specExp = 20.0;
half3 specularCol = h3tex2D( g_eyeSpecMap, v2f.UVtex0 ) * g_specAmount;

float tea;

float3 centerToSurfaceVec = objFlatNormal; // = normalize( v2f.objCoord )
float firstDot = centerToSurfaceVec.y; // = dot( ce
if( firstDot > 0.805000 )
{
  // We hit the iris.  Do the math.

  // we start with a ray from the eye to the surface
  float3 ray_dir = normalize( v2f.objCoord - objEye
  float3 ray_origin = v2f.objCoord;

  // refract the ray before intersecting with the iris
  ray_dir = refract( ray_dir, objFlatNormal, g_refra

  // first, see if the refracted ray would leave the e
  float t_eyeballSurface = SphereIntersect( 16.0, r
  float3 objPosOnEyeBall = ray_origin + t_eyeball
  float3 centerToSurface2 = normalize( objPosOn

  if( centerToSurface2.y > 0.805000 )
  {
    // Display a blue color
    diffuseCol = float3( 0, 0, 0.7 );
    objBumpNormal = objFlatNormal;
    specularCol = float3( 0, 0, 0 );
    specExp = 10.0;
  }
  else
  {
    // transform into irisSphere space
    ray_origin.y -= 5.109000;

    // intersect with the Iris sphere
    float t = SphereIntersect( 9.650000, ray_origin, ray_dir );
    float3 SphereSpaceIntersectCoord = ray_origin + t * ray_dir;
    float3 irisNormal = normalize( -SphereSpaceIntersectCoord );
```

Eye Shader from Luna
Maximum branch takes 674 cycles
Minimum branch takes 193 cycles.



```
C:\WINDOWS\System32\cmd.exe

T:\tmp>t:\sw\devrel\sdk\tools\bin\release_pdb\nvshperf\nvshaderperf -a NV40 corn
ea2.txt
------------------------------------------------------------------------------

Running performance on file Cornea2.txt
------------------------ NV40 ------------------------
Target: GeForce 6800 Ultra (NV40) :: Unified Compiler: v77.72
Cycles: 674.25 :: R Regs Used: 12 :: R Regs Max Index (0 based): 11
Pixel throughput (assuming 1 cycle texture lookup) 9.50 MP/s

T:\tmp>t:\sw\devrel\sdk\tools\bin\release_pdb\nvshperf\nvshaderperf -minbranch -
a NV40 cornea2.txt
------------------------------------------------------------------------------

Running performance on file Cornea2.txt
------------------------ NV40 ------------------------
Target: GeForce 6800 Ultra (NV40) :: Unified Compiler: v77.72
Cycles: 192.82 :: R Regs Used: 12 :: R Regs Max Index (0 based): 11
Pixel throughput (assuming 1 cycle texture lookup) 33.33 MP/s

T:\tmp>_
```

# NVShaderPerf – version 2.0

- Vertex throughput
- GLSL vertex program
- Multiple driver versions from one NVShaderPerf
- What else do you need?

[NVShaderPerf@nvidia.com](mailto:NVShaderPerf@nvidia.com)

# NVPerfKit

- **What is NVPerfKit?**
- **Associated Tools**
- **NVPerfKit 2.0**

# NVPerfKit: The Solution!

- Why is my app running at 13FPS after CPU tuning?
- How can I determine what is going in that GPU?
- How come IHV engineers are able to figure it out?

# What is NVPerfKit?

- **Driver and GPU performance counters**
    - **Performance Data Helper (PDH)**
    - **Microsoft PIX for Windows**
- **NVPerfHUD functionality inside any application**
- **Application triggered sampling**
- **OpenGL and Direct3D**

# NVPerfKit: What it looks like...

# What is in the NVPerfKit package?

- **Instrumented Driver**
  - **Exposes GPU and Driver Performance Counters**
  - **Supports OpenGL and Direct3D**
  - **Supports SLI Counters**
- **Tools**
  - **NVDevCPL**
  - **PIX Plugin**
  - **NVAppAuth**
- **SDK**
  - **Sample Code**
  - **Helper Classes**
  - **Docs**

# OpenGL Signals

| Counter Description | Official Name |
|---|---|
| FPS | OGL FPS |
| Frame Time (1/FPS) | OGL frame time mSec |
| Driver Sleep Time (driver waits for GPU) | OGL frame mSec Sleeping |

# Direct3D Signals

| Counter Description | Official Name |
|---|---|
| FPS | D3D frame FPS |
| Frame Time (1/FPS) | D3D frame time mSec |
| AGP Memory Used | D3D frame agpmem MB |
| Video Memory Used | D3D frame vidmem MB |
| Driver Time | D3D frame mSec in driver |
| Driver Sleep Time (driver waits for GPU) | D3D frame mSec Sleeping |
| Triangle Count | D3D frame tris |
| Batch Count | D3D frame num batches |
| Locked Render Targets Count | D3D Locked Render Targets |

# GPU Signals

**GPU**

**Vertex Setup**

**Vertex Shader**

**Rasterizer**

**Texture**

**Pixel Shader**

**Frame Buffer**

**gpu_idle**
**vertex_attribute_count**

**vertex_shader_busy**

**culled_primitive_count**
**primitive_count**
**triangle_count**
**vertex_count**

**fast_z_count**
**shaded_pixel_count**

**shader_waits_for_texture**

**pixel_shader_busy**

**shader_waits_for_rop**

Supported GPUs
Quadro FX 4500
GeForce 7800 GTX
GeForce 6800 Ultra & GT
GeForce 6600

# NVPerfKit Demo: Pixel Shader Bound

# NVPerfKit Demo: Texture Bound

# What is PDH? What is Perfmon?

- **PDH: Performance Data Helper for Windows**
  - **Win32 API for exposing performance data to user applications**
  - **Standardized interface with many providers and clients**

- **Perfmon: (aka Microsoft Management Console)**
  - **Win32 PDH client application**
  - **Low frequency sampling (1X/s)**
  - **Displays PDH based counter values:**
    - **OS: CPU usage, memory usage, swap file usage, network stats, etc.**
    - **NVIDIA: all of the signals exported by NVPerfKit**

# Associated Tools: Perfmon

# Associated Tools: NVDevCPL

# Associated Tools: NVIDIA Plug-In for Microsoft PIX for Windows

# Associated Tools: NVIDIA Plug-In for Microsoft PIX for Windows

# Helper Classes and Code Samples

- **CPDHHelper: simplifies using PDH**

```
int nIndex = pdh.add("countername");
pdh.sample();
float fValue = pdh.value(nIndex);
```

- **CTrace: ring buffer for holding performance data**
- **CTraceDisplay: simple API agnostic graphing library**
- **OpenGL and Direct3D sample apps**
  - Integration of helper classes
  - Security mechanism usage

# Graphic Remedy's gDEBugger 2.0

# NVPerfKit 2.0

- **Simplified Experiments**
- **Targeted analysis to ease bottleneck determination**
  - **Supplement PDH based single counters**
  - **Multi-pass experiments where multiple GPU counters are needed to compute results**
  - **Exposes all of the power of NVPerfHUD 4.0 to developers**
- **More OpenGL and Direct3D counters**
- **NVPerfHUD 4.0**
- **Linux support**

# Simplified Experiments

- Usage:

```
NVPMAddCounter("ps_utilization");
NVPMAddCounter("vs_utilization");
NVPMAddCounter("gpu_idle");
NVPMAllocObjects(50);

int nNumPasses;
NVPMBeginExperiment(&nNumPasses);
for(int ii = 0; ii < nNumPasses; ++ii) {
    NVPMBeginPass(ii);

    // Draw the frame
    NVPMBeginObject(0);
    // DPs associated with object 0
    NVPMEndObject(0);

    NVPMBeginObject(1);
    // DPs associated with object 1
    NVPMEndObject(1);

    // ...
    NVPMEndPass(ii);
}
NVPMEndExperiment();
NVPMGetCounterValue("ps_utilization", 0, &fPSSol); // 0 == object id
NVPMGetCounterValue("vs_utilization", 0, &fVSSol);
```

**NVPerfHUD 4.0**

Raul Aguaviva

# Agenda

- **What is NVPerfHUD?**
- **How does it work?**
- **Demo**
- **Schedule**

# What is NVPerfHUD?

- **Stands for: PERFormance Heads Up Display**
  - **Overlays graphs and dialogs on top of your application**
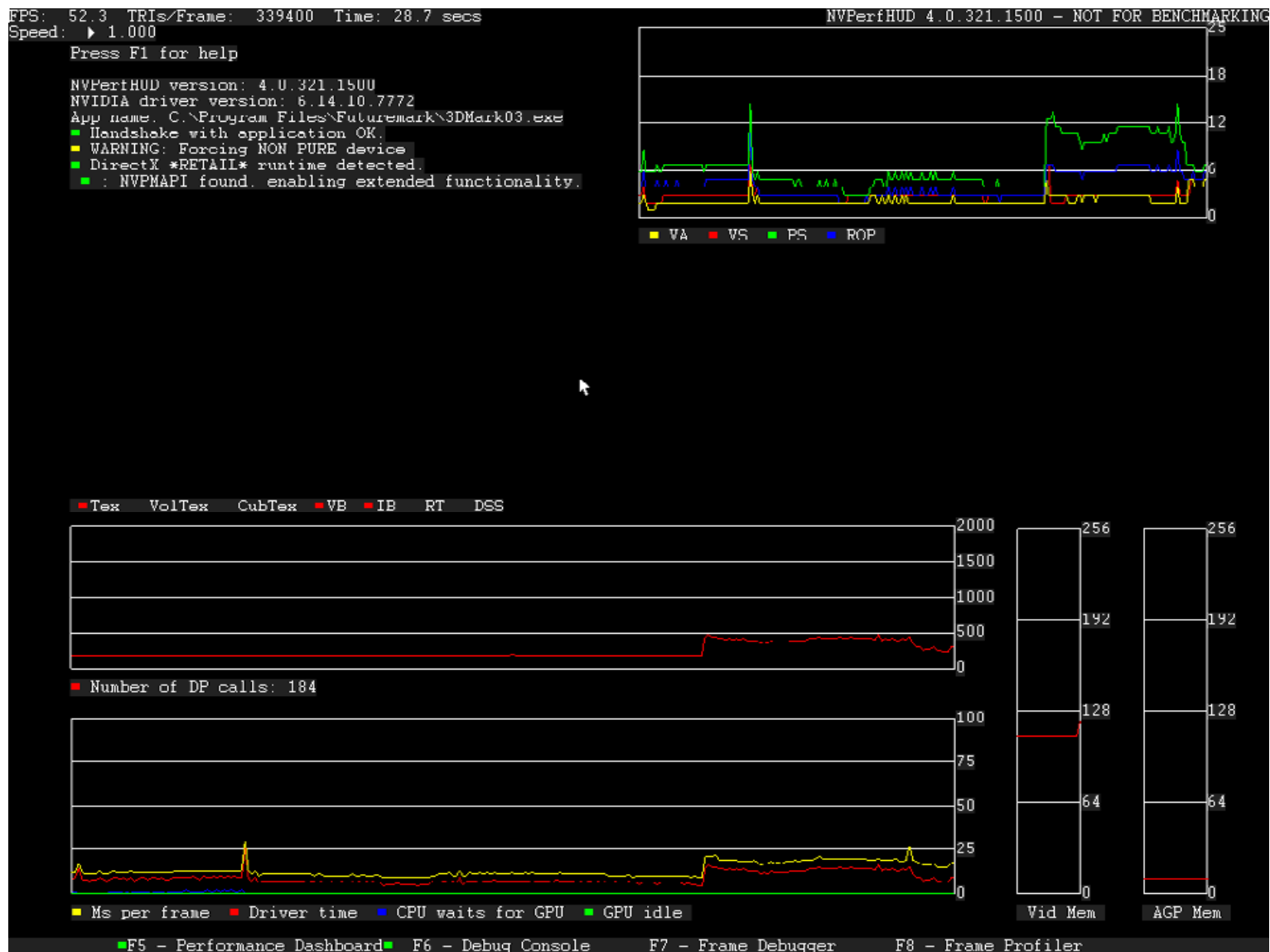  - **Interactive HUD**

# What is NVPerfHUD?

- **4 different types of HUD**
  - **Performance Dashboard**
  - **Debug Console**
  - **Frame Debugger**
  - **Frame Profiler  (New in 4.0)**

# How to use it

- **Run your application with NVPerfHUD**
- **Use it as you normally do until you find:**
  - **Functional problem: use the debugger**
  - **Low FPS: use the profiler**

# Performance Dashboard
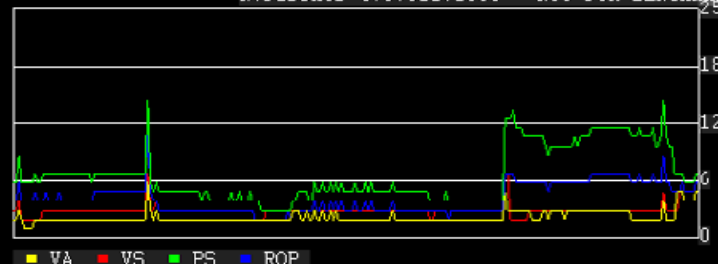
# Performance Dashboard

# Performance Dashboard

# Performance Dashboard

# Performance Dashboard

# Performance Dashboard

# Performance Dashboard

- **Resource monitor**

```
 Tex    VolTex    CubTex   VB   IB    RT    DSS
```

- **Resources monitored**
  - Textures
  - Volume Textures
  - Cube textures
  - Vertex Buffers
  - Index buffers
  - Stencil and depth surfaces

# Performance Dashboard

# Performance Dashboard

- **Speed control**

```
FPS:    52.3   Tris/Frame:    339400   Time: 28.7
Speed:      ▶  1.000
          Press F1 for help

          NVPerfHUD version: 4.0.321.1500
          NVIDIA driver version: 6.14.10.7772
          App name: C:\Program Files\Futuremark\
```

# The simplified graphics pipeline

# Performance Dashboard Demo

- **Install**
- **Configure**
- **Drag & Drop**

# Debug Console

# Frame Debugger

- **Demo**

# Frame Debugger, advanced view

- **Demo**

# Frame Profiler

- **Measures performance counters**
- **strategy**

# Frame Profiler, measuring

- **NVPerfHUD uses NVPerfKit**
  - **uses ~40 Performance Counters (PC's)**
- **Can not read all of them at the same time**
- **Need to render THE SAME FRAME until all the PC's are read**

# Frame Profiler, strategy

- **Optimization Strategy:**
  - **Group by state is roughly grouping by bottleneck**
  - **These groups are called "state buckets"**

- **Procedure**
  - **Group draw calls by rendering state into state buckets**
  - **Identify the bottleneck of the most expensive state bucket**
    - **Solved by NVPerfHUD**
  - **Cure the bottleneck with a common corrective action**

# Frame Profiler Demo

# Frame Profiler Demo, advanced view

# About freezing the application

- Only possible if the application uses time-based animation

- Stop the clock
  - Intercept: QueryPerformanceCounter(), timeGetTime()
  - NO RDTSC!!

- Pos += V * DeltaTime

# Schedule

- **Beta: August**
- **Release : September**

NVPerfHUD 3 - Mozilla Firefox

File  Edit  View  Go  Bookmarks  Tools  Help

http://developer.nvidia.com/object/nvperfhud_home.html    Go    nvperfhud 3.0

Firefox Help    Firefox Support    Plug-in FAQ    Sign In - Yahoo!    Uber Gizmo, the Ga...    Microsoft Visual C++...    DirectX    Amazon.com: All Pr...    sapere.it

developer.nvidia.com
THE SOURCE FOR GPU PROGRAMMING

Search
www   developer.nvidia.com

nVIDIA

REGISTERED
DEVELOPERS
  JOIN

HOME
NEWS
DOCUMENTATION
TOOLS & SDKs
PARTNERS

NEWSLETTER SIGN-UP
EVENTS CALENDAR
DRIVERS
CONTACT

Legal Info

## NVPerfHUD 3

### Quick Links

+ Introductory Video
+ Downloads
+ GameDev.net Review

### Overview

Modern GPUs generate images through a pipelined sequence of operations.  A pipeline runs only as fast as its slowest stage, so tuning graphical applications for optimal performance requires a pipeline-based approach to performance analysis. NVPerfHUD analyzes your graphics pipeline performance and provides real-time statistics you can use to diagnose performance bottlenecks in your Direct3D application.

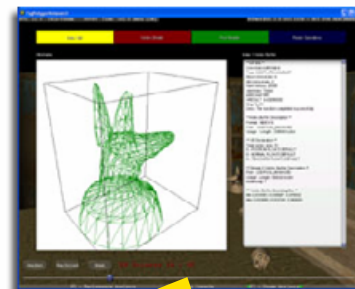The latest release includes several new features and separate modes:

+ **Frame Analysis Mode**
  Freeze your application and single-step through the current frame. You can see what is happening inside the GPU at each stage of your graphics pipeline including the Index Unit, Vertex Shader, Pixel Shader & Render Operations Unit.

+ **Debug Console Mode**
  This mode shows you DirectX Debug Runtime messages, NVPerfHUD, and custom messages from your application.

+ **Performance** 
  The powerful performance analysis experiments from the previous release are still available in Per...

**Now Available!**

The opt-in mechanism... no code alterations are necessary if you have already enabled NVPerfHUD 2.0 in your...

Be sure to check out the Getting Started instructions in the **NVPerfHUD User Guide**, and read through the methodology for effectively identifying and crushing performance bottlenecks in your application. We've also created a **Quick Reference Card** with tips and shortcuts that you can keep at your fingertips. Both these documents are available in English, Japanese, Chinese, and Korean.

See our "**NVIDIA Performance Analysis Tools**" talk from **GDC 2005** for more information on how to analyze your applications using NVPerfHUD. In addition, our "**Practical Performance Analysis and Tuning**" talk from **GDC 2004** explains the theory of pipeline analysis and bottleneck removal.

### Downloads

NVPerfHUD

Quick Reference

## Developer Reviews

NVPerfHUD has way more features than I expected! Blew me away!!
  - *Remy Saville*
    *Graphics Programmer*
    *Relic Entertainment, Inc.*

[Frame Analysis Mode] is by far the most impressive mode of NVPerfHUD. ... A few moments in this mode does more to show how advanced graphic pipelines work than anything else I have seen. ... NVPerfHUD should be of interest to anybody who develops software that employs DirectX graphics. Even if you don't have an NVIDIA card it's worth the price of a GeForce 6600 or 6800 to get this tool.
  - *Bryan Mau*
    *GameDev.net Review*

NVPerfHUD 3 - it's simply amazing.  I use it virtually every day.
  - *Chris King, President*
    *IDV, Inc.*

NVPerfHUD is a great tool for debugging and performance analysis.
  - *Richard Schubert*
    *Graphics/Effects Programmer*
    *Yager Development GmbH*

I just wanted to drop a note to thank you for the hard work on the NVPerfHUD. This new version is incredibly useful. Not a day goes by that I don't appreciate the support NVIDIA gives developers with tools such as this.
  - *Matt Shaw*
    *Director of Technology*
    *Mythic Entertainment*

Done

Start    2 Win...    Oxford ...    4 Fire...    4 Micr...    Total C...    nvPerf...    unan^M...    Microso...    5:59 PM

# Questions?

- **Developer tools DVDs available at our booth**
- **Online: http://developer.nvidia.com**

NVGLExpert@nvidia.com

NVShaderPerf@nvidia.com

NVPerfKIT@nvidia.com

NVPerfHUD@nvidia.com

FXComposer@nvidia.com

# SLI

**Matthias M Wloka**

**NVIDIA Corporation**

# SLI: Scalable Link Interface

- **Plug 2 identical GPUs into PCI-E motherboard**

- **Driver still reports only one (logical) device**
  - **Renders up to 1.9x faster**

- **Video memory does NOT d**

# Don't Care For High-End Niche Markets

- **SLI becoming mainstream:**
  - **GeForce 6600 GT SLI**
  - **In addition to 6800 GT and 6800 Ultra**

- **Dual core boards**
  - **Gigabyte 3D1: Dual 6600 GT**

- **SLI motherboards sold to date: > 350,000 units**
  - **That's > 25% of total nForce 4**

# Game Development Cycle

- **2 years (or more)**
  - **CPU performance doubles (or less)**
  - **GPU performance quadruples**

- **CPU/GPU balance shifts!**
  - **Worse: CPU-hungry modules come later: AI, physics, full game play**

- **SLI hints at future GPU vs. CPU balance**
  - **For target 'mainstream' spec**

# The Last Couple of Years



**Courtesy Ian Buck, Stanford University**

# Ok, How Does SLI Work?

- **Compatibility mode:**
  - **Only uses one GPU**
  - **No SLI benefits**

- **Alternate frame rendering (AFR)**

- **Split frame rendering (SFR)**

# AFR

- **Each GPU works on its own frame**

GPU 0:  **1**  **3**  ...

GPU 1:  **2**  **4**  ...

- **Scan-out toggles where to read framebuffer from**

# General Rendering Case for AFR

- **If frame not self-contained:**
  - **Push necessary data to other GPU**
  - **E.g., updating render-to-texture targets only every other frame**

- **Pushing data to other GPU is overhead**
  - **Hence not 2x speed-up**

# AFR Advantages

- **All work is parallelized**
  - **Pixel fill, raster, vertex transform**

- **Preferred SLI mode**

- **Works best when frame self-contained**
  - **No prior work is re-used**
  - **No communications overhead between GPUs**

# SFR

- **Both GPUs work on the same frame**
  - GPU 0 renders top portion
  - GPU 1 renders bottom portion

| GPU 0 | | | |
|---|---|---|---|
| GPU 1 | 1 | 2 | 3 | ...

- **Scan-out combines framebuffer data**

# General Rendering Case for SFR

- **Load-balance 'top' vs. 'bottom'**
  - **If one GPU took longer to render**
  - **Adjust load accordingly (make it work less)**

- **Clip vertices to top/bottom portions**
  - **Avoids both GPUs processing all vertices**
  - **But not perfect**

- **Still requires data sharing:**
  - **E.g., render to texture**

# SFR Compared to AFR

- SFR works even when limiting number of frames buffered
  - Or when AFR otherwise fails

- In general, SFR has more communications overhead

- Applications with heavy vertex load benefit less from SFR

# How Do I Detect SLI Systems?

- **NVCpl API:**
  - NVIDIA-specific API supported by all NV drivers

- **Function support for:**
  - Detecting that NVCpl API is available
  - Bus mode (PCI/AGP/PCI-E) and rate (1x-8x)
  - Video RAM size
  - SLI

# NVCpl API SLI Detection

- SDK sample and full documentation available

```
HINSTANCE hLib = ::LoadLibrary("NVCPL.dll");

NvCplGetDataIntType NvCplGetDataInt;
NvCplGetDataInt =
    (NvCplGetDataIntType)::GetProcAddress(hLib,
                        "NvCplGetDataInt");


long     numSLIGPUs  = 0L;
NvCplGetDataInt(NVCPL_API_NUMBER_OF_SLI_GPUS,
                        &numSLIGPUs);
```

# Forcing SLI Support In Your Game

- **Use NVCpl**
  - **NvCplSetDataInt() sets AFR, SFR, Compatibility mode**
  - **See SDK sample**

- **Modify or create a profile:**
  - **http://nzone.com/object/nzone_sli_appprofile.html**
  - **End-users can create profiles as well**

# Overview: Things Interfering with SLI

- CPU-bound applications
  - Or vsync enabled

- Limiting number of frames buffered

- Communications overhead

# CPU-Bound Applications

- SLI cannot help

- Reduce CPU work or better:

- Move CPU work onto the GPU
  - See http://GPGPU.org

- Don't throttle frame-rate

# VSync Enabled

- **Throttles frame-rate to monitor refresh**

- **Enabling triple-buffering does NOT offset enabling vsync:**
  - **If render-rate faster than monitor refresh,**
  - **Then vsync still gates GPU**

- **Worse, triple-buffering**
  - **Increases lag**
  - **Consumes (much) more video-memory**

# Limiting Number of Frames Buffered

- **Some apps allow at most one frame buffered**
  - **To reduce lag**
  - **Via event queries**
  - **Don't lock/read back-buffer: Causes CPU stall!**

- **Disables AFR SLI speed-up**

- **But SLI is up to ~1.9x faster**
  - **I.e., SLI systems ~1.9x less lag**

# Why Locking the Back-Buffer Is Bad

Back-buffer lock:

wait for GPU to finish rendering

| CPU | | CPU | |
|---|---|---|---|

GPU | | GPU | |

Frame n          Frame n+1                    ...

# Limit Frames Buffered to Number of GPUs

- **Single GPU system:
  Buffer at most 1 frame**

- **When detecting SLI system:
  Buffer at most 2 frame**

# The Basic Pipeline

CPU → Push Buffer → GPU0 →

## Frames flow through pipe over time:

| | | | |
|---|---|---|---|
| Frame n | CPU | Push Buffer | GPU0 |
| Frame n+1 | CPU | Push Buffer | GPU0 |
| Frame n+2 | CPU | Push Buffer | GPU0 |
| Frame n+3 | CPU | Push Buffer | GPU0 |

1 Frame = L ms

Time

# Single GPU Latency

| | | |
|---|---|---|
| CPU | Push Buffer | GPU0 |

User inputs data

| | | | |
|---|---|---|---|
| CPU | Push Buffer | GPU0 | } 1 Frame = L ms |

CPU processes it

PB buffers it

| | | |
|---|---|---|
| CPU | Push Buffer | GPU0 |

GPU processes it
Result visible

| | | |
|---|---|---|
| CPU | Push Buffer | GPU0 |

Time

**Total latency: 3L ms**

# Latency Assumptions

- **GPU limited**
  - If not, then push buffer contains <1 frame
  - No point in limiting push buffer

- **SLI is 2x faster**
  - Can relax this later!

- **Increase frames buffered to 2:**

| → | CPU | Push Buffer | Push Buffer | GPU0 | GPU1 | → |

# Frames Flowing Through AFR SLI

| | | | | |
|---|---|---|---|---|
| CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |

**Frame n**

| | | | | |
|---|---|---|---|---|
| CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |

} 1 Frame = L/2 ms

**Frame n+1**

| | | | | |
|---|---|---|---|---|
| CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |

**Frame n+2**

| | | | | |
|---|---|---|---|---|
| CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |

**Frame n+3**

| | | | | |
|---|---|---|---|---|
| CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |

**Frame n+4**

| | | | | |
|---|---|---|---|---|
| CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |

**Frame n+5**

**Time**

# AFR SLI Latency

| | CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |
|---|---|---|---|---|---|
| | CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |
| | CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |
| | CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |
| | CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |
| | CPU | Push Buffer | Push Buffer | GPU0 | GPU1 |

User inputs data
CPU processes it

} 1 Frame = $L/2$ ms

PB buffers it

PB buffers it

GPU processes it

GPU processes it
Result visible

Time

**Total latency: $5 \cdot L/2$ ms**

# Latency Comparison: Single vs. AFR

- **Single GPU latency: 3L ms**
  - **3 frames of length L ms**

- **AFR SLI GPU latency: 5 L/2 = 2.5L ms!**
  - **5 frames of length L/2 ms (i.e., double frame rate)**
  - **Despite buffering twice as many frames**

- **SLI speed-up only needs to be 1.66!**
  - **$3L = 5L/x \rightarrow x = 5L/3L = 1.66$**
  - **Most games speed-up by ~1.8**

# SFR Latency?

- **SFR unaffected by buffering one frame**

- **SFR speed-up directly reduces lag**
  - **If SFR 2x faster,**
  - **Then latency 2x shorter**

# Even Better: Limit Lag Based on FPS

- **If your game runs at over 100 fps**
  - **Reasonable to buffer 3 frames**

- **If your game runs at less than 15 fps**
  - **Only allow one frame to buffer**

- **Faster SLI system gets automatic benefit**

- **Our drivers already do that**
  - **> 15 fps buffer 3 frames as usual**
  - **< 15fps reduce number of frames buffered**

# Overview: Things Interfering with SLI

- CPU-bound applications
  - Or vsync enabled

- Limiting number of frames buffered

- **Communications overhead**

# Communications Overhead

- Peer to peer SLI memory transfers
  - Transfer itself costs bandwidth and time
  - GPU stalls waiting for transfer to complete

- Or replicate operations on both GPUs
  - For example, render to texture

- Relevant resources:
  - Vertex/index buffers
  - Textures
  - Render targets

# Uploading Resources On the Fly

- Remember video RAM is duplicated

- Need to transfer to both video RAMs

- Not much developers can do to avoid this
  - Oh well

# Render Targets

- **Clear Z**
  - **Always clear Z!**

- **Clear color when detecting SLI**
  - **Tells driver that the old data is irrelevant**
  - **No need to transfer old data across GPUs**

- **Don't reuse data across frames**
  - **Make frames self sufficient, i.e., independent from one another**

# Update-Skipping "Optimization"

| Frame n | Frame n+1 | Frame n+2 | Frame n+3 |
|---------|-----------|-----------|-----------|
| RTT   Use T | Use T | RTT   Use T | Use T |

↑ **Needs RTT**  (under Frame n+1)

↑ **Needs RTT**  (under Frame n+3)

- **Added SLI overhead:**

| GPU 0 | GPU 1 | GPU 0 | GPU 1 |
|-------|-------|-------|-------|

- **GPU 1 stalls until GPU 0 RTT finishes and transfers**
- **Or GPU 1 duplicates RTT operation**
- **Might as well do right thing when on SLI**
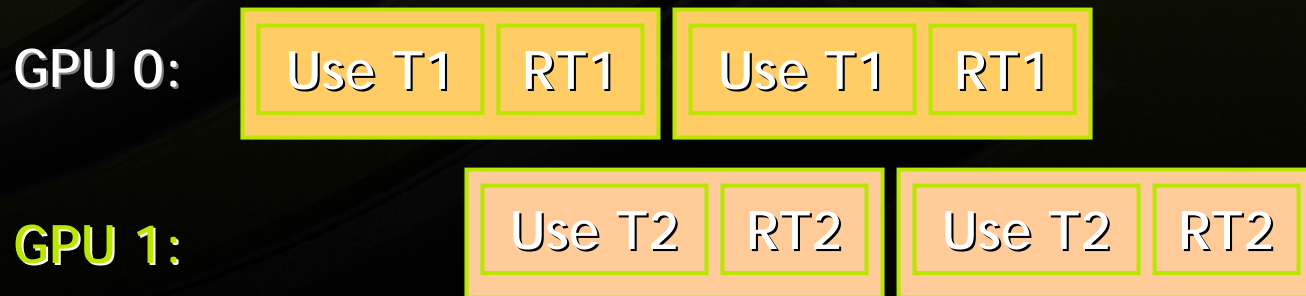
# Render Early, Use Late!

GPU 0:

| RT1 | Use T1 | RT1 | Use T1 |

GPU 1:

| | Use T1 | | Use T1 |

- **Avoid sync-stalls**
  - **In AFR SLI as shown**
  - **And in single GPU mode**
  - **But still has communications overhead**

# Really Bad: Use Early, Render Late

**GPU 0:** | Use T1 | RT1 |     | Use T1 | RT1 |

**GPU 1:** | Use T1 | RT1 |     | Use T1 | RT1 |

## Instead: Ring-buffer textures when on SLI!

**GPU 0:** | Use T1 | RT1 | | Use T1 | RT1 |

**GPU 1:** | Use T2 | RT2 | | Use T2 | RT2 |

# SLI Performance Debug Support

- ## SLI support in NVPerfKit:
  - Pluggable hardware and driver signals for
  - PIX
  - perfmon.exe
  - pdh (your game, VTune…)

- ## "NVIDIA Performance Analysis Tools" Today, 2:30pm - 3:30pm

# Supported SLI Performance Signals

- Total SLI peer-to-peer bytes
- Total SLI peer-to-peer transactions

- Above originating from
  - Vertex/index buffers: bytes and transactions
  - Textures: bytes and transactions
  - Render targets: bytes and transactions

# Questions?

- **GPU Programming Guide, Chapter 8**
  **http://developer.nvidia.com/object/gpu_programming_guide.html**

- **http://developer.nvidia.com**
  **The Source for GPU Programming**

- **mwloka@nvidia.com**

- **Slides available online**

# NVIDIA SDK
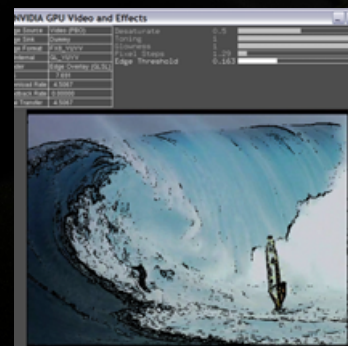
## *The Source for GPU Programming*



Hundreds of code samples and effects that help you take advantage of the latest in graphics technology.

- Tons of updated and all-new DirectX and OpenGL code samples with full source code and helpful whitepapers:

  **Transparency AA, GPU Cloth, Geometry Instancing, Rainbow Fogbow, 2xFP16 HRD, Perspective Shadow Maps, Texture Atlas Utility, ...**

- Hundreds of effects, complete with custom geometry, animation and more:
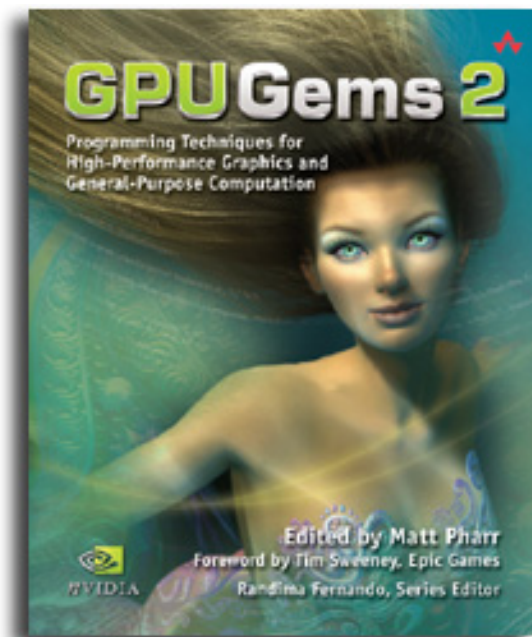
  **Shadows, PCSS, Skin, Plastics, Flame/Fire, Glow, Image Filters, HLSL Debugging Techniques, Texture BRDFs, Texture Displacements, HDR Tonemapping, and even a simple Ray Tracer!**
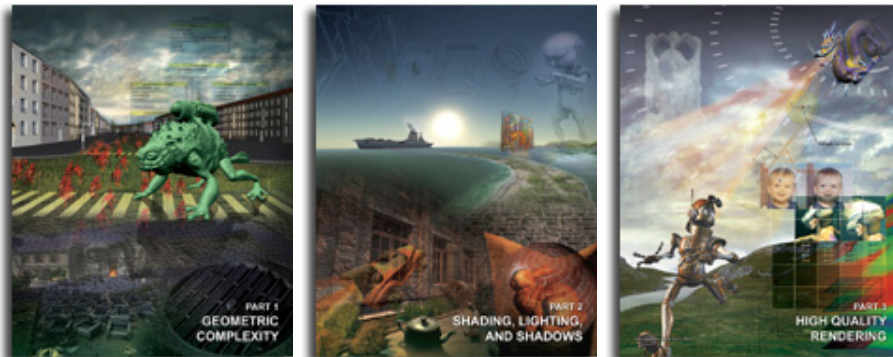
# GPU Gems 2

## Programming Techniques for High-Performance Graphics and General-Purpose Computation

- **880 full-color pages**
- **330 figures**
- **Hard cover**
- **$59.99**
- **Experts from universities and industry**

## Graphics Programming

- Geometric Complexity
- Shading, Lighting, and Shadows
- High-Quality Rendering

## GPGPU Programming

- General Purpose Computation on GPUs: A Primer
- Image-Oriented Computing
- Simulation and Numerical Algorithms