# GPU Performance Tools and Analysis Techniques

## European Developer's Forum 2004

### Sébastien Dominé

### NVIDIA

# Agenda

- Performance Tools Survey

- Practice

- Next generation Performance Tools

- Conclusion

- Q & A

# Agenda

- Performance Tools Survey

- Practice

- Next generation Performance Tools

- Conclusion

- Q & A

# Performance Tools Survey

- NVPerfHUD
  - Direct3D9 Performance HUD
- NVShaderPerf
  - Offline Shader Performance Analysis
- FX Composer
  - HLSL Shader Editor IDE

# NVPerfHUD 2.0

- Overlay graph that displays stats from :
  - Direct3D9 API interception layer
  - Direct3D Driver
  - Requires NVIDIA GPU
- Able to bypass and inject API calls to assist with performance analysis
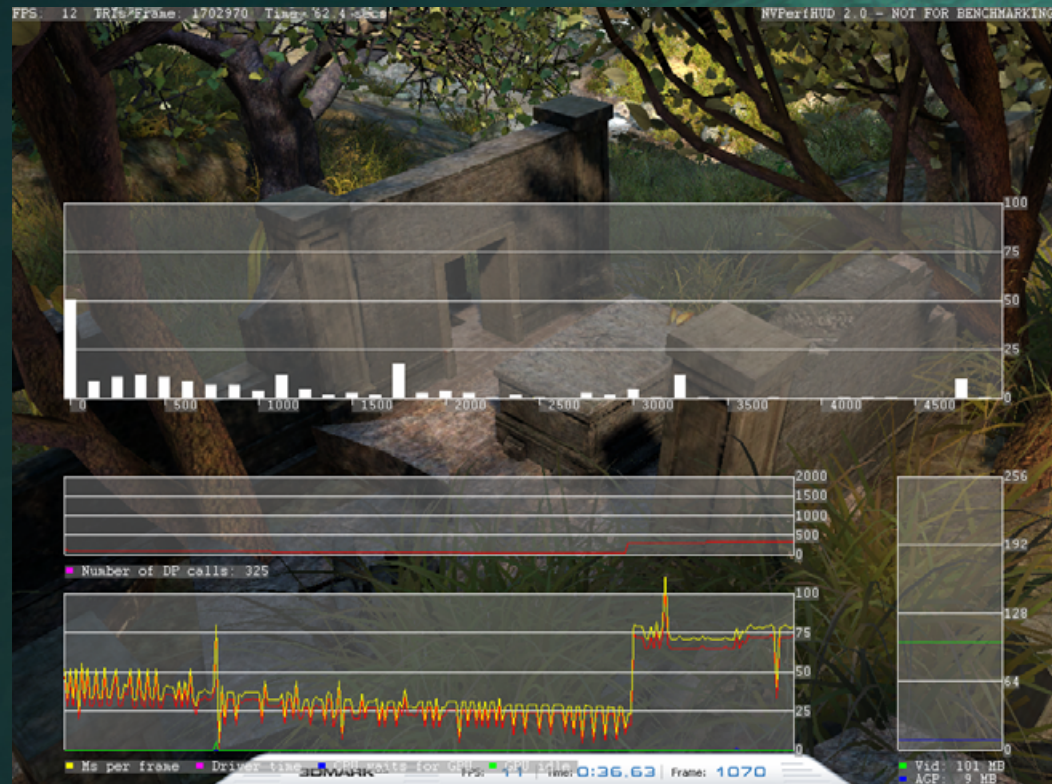- Only works on your own application



Image courtesy of FutureMark Corp.

FPS: 12  TRIs/Frame: 1702970  Time: 62.4 secs

NVPerfHUD 2.0 — NOT FOR BENCHMARKING

- Frame rate
- Number of triangles/frame
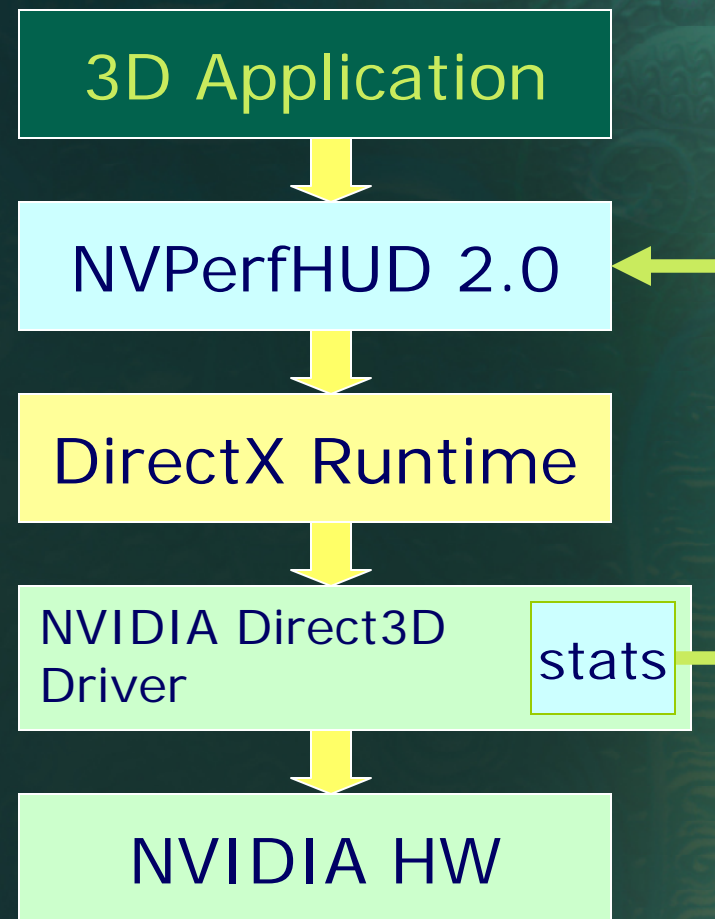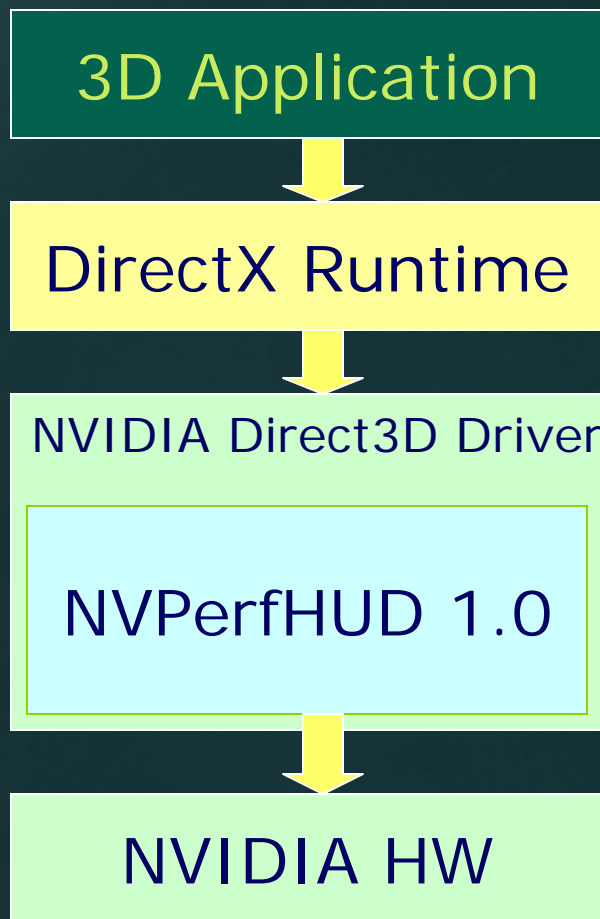- Elapsed time in the session

es Batches:

- Driver waiting for GPU (Spin)
- GPU Idle Performance Counter

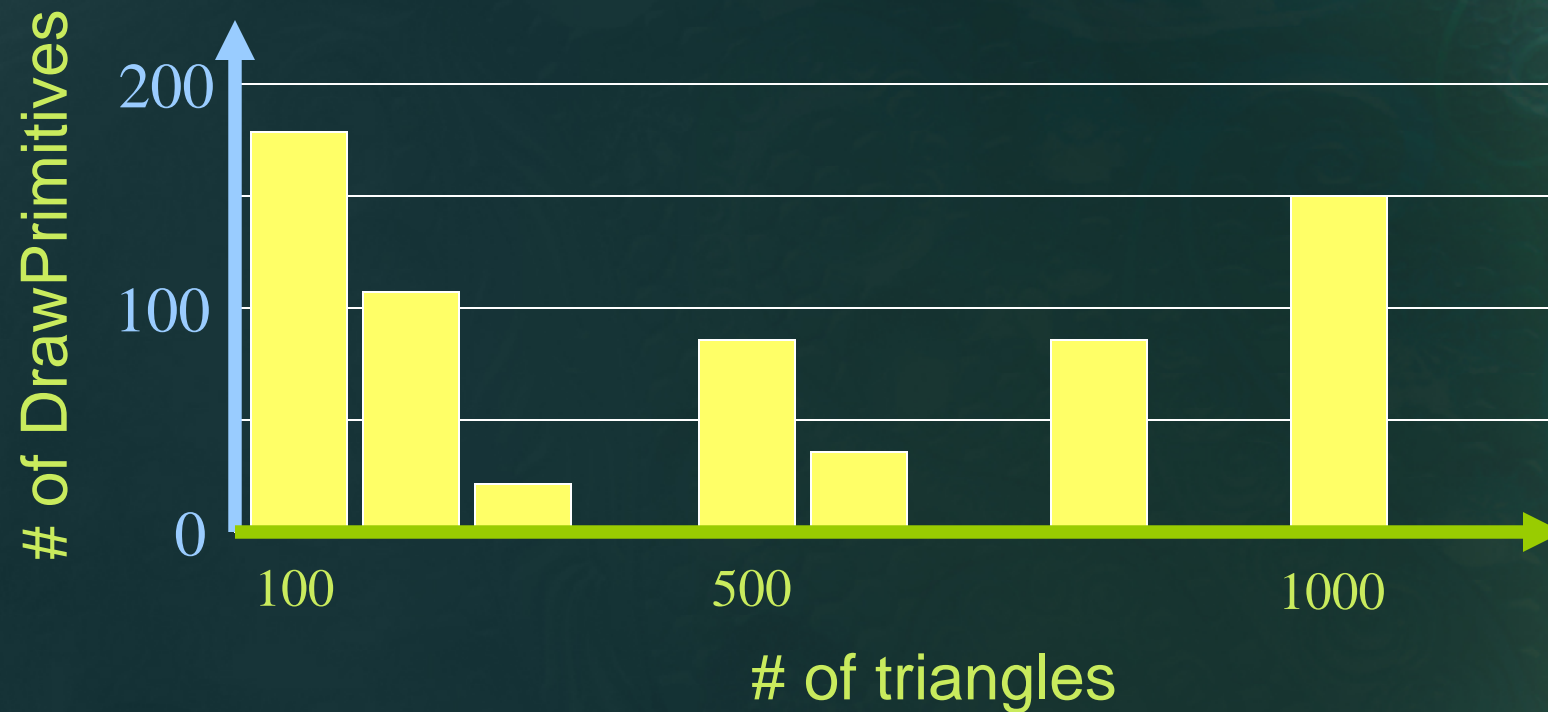Current Memory footprint:

Histogram of Draw Primitives Batches

Video Memory

Number of

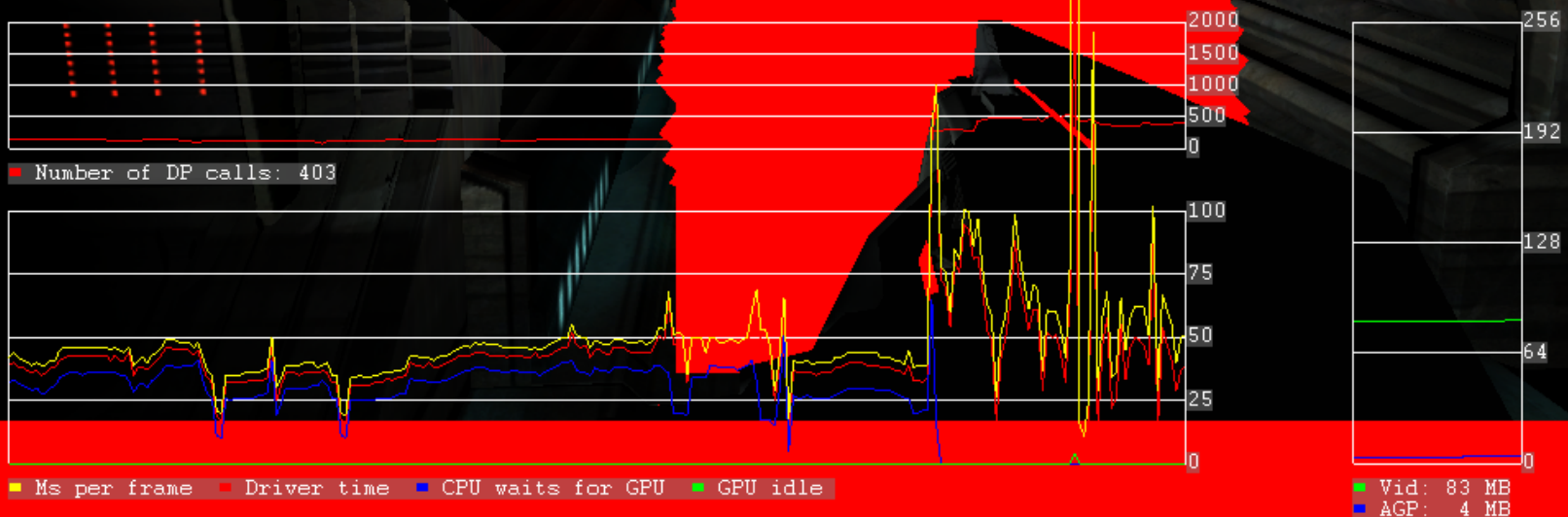Ms per frame   Driver time   CPU waits for GPU   GPU idle

3DMARK   FPS: 11   Time: 0:36.63   Frame: 1070

Vid: 101 MB
AGP:   9 MB

Image courtesy of FutureMark Corp.

# What's new in 2.0?

**Left diagram:**

3D Application
↓
DirectX Runtime
↓
NVIDIA Direct3D Driver
- NVPerfHUD 1.0
↓
NVIDIA HW

→

**Right diagram:**

3D Application
↓
NVPerfHUD 2.0
↓
DirectX Runtime
↓
NVIDIA Direct3D Driver → stats
↓
NVIDIA HW

- DrawPrimitives/DrawIndexedPrimitives Histogram

Texture Stage States

Image courtesy of FutureMark Corp.

Pixel Shaders 1.x

Image courtesy of FutureMark Corp.

Pixel Shaders 2.0
Pixel Shaders 3.0

Image courtesy of FutureMark Corp.

# 2x2 Texture replacement

FPS: 19   TRIs/Frame:   344384   Time: 78.2 secs

NVPerfHUD 2.0 — NOT FOR BENCHMARKING

Number of DP calls: 157

Ms per frame   Driver time   CPU waits for GPU   GPU idle

Vid: 79 MB
AGP: 15 MB

Image courtesy of FutureMark Corp.

Image courtesy of FutureMark Corp.

# NVPerfHUD - Overhead

- NVPerfHUD is fairly lean but...
- Overlay graph and DLL interception can costs up to 1.3%
- Driver instrumentation can cost up to 6%
- Upper bound for total cost: 7%

# NVShaderPerf

Inputs:
- HLSL
- PS1.x, PS2.x, PS3.x
- VS1.x, VS2.x, VS3.x

GPU Arch:
- GeForce FX (NV3X)
- GeForce 6 Series (NV4X)
- Quadro FX (NV3X+NV4X)

**NVShaderPerf**

Outputs:
- Assembly code
- # of cycles
- # of temporary registers
- Pixel throughput
- Forces all fp16 and all fp32 (gives performance bounds)



```
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 r1.w, r2.w, r2.w
mul r1.w, r1.w, r2.w
mov r2.x, c9.w
add r2.x, r2.x, -c8.x
mad r1.w, r1.w, r2.w, c8.x
dp3 r0.x, r0, r1
mul r0.w, r0.w, r1.w
mul r1.xyz, r0.w, c4
```
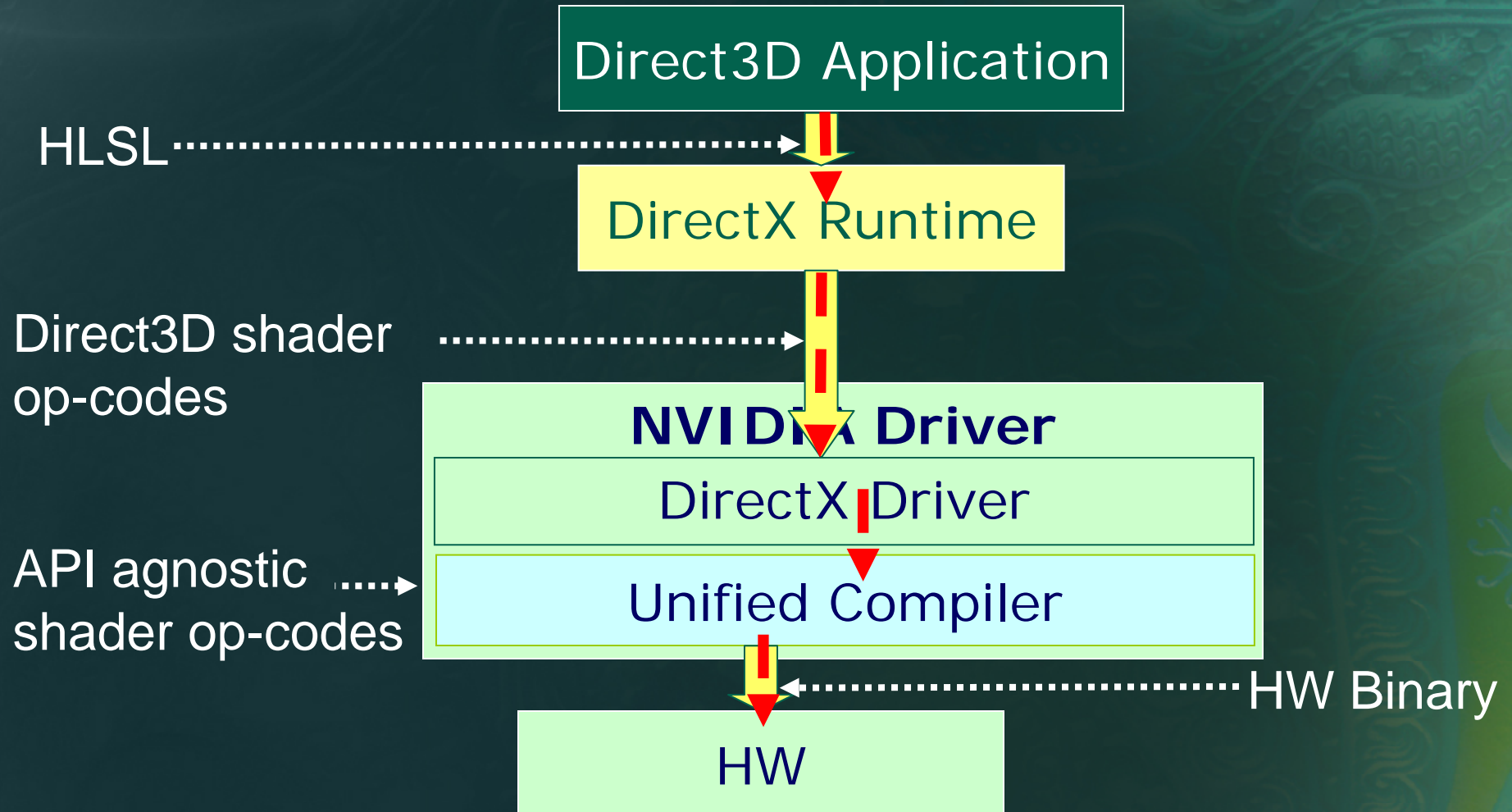
```
Shader performance using all fp32
Cycles: 21.00 :: R Regs Used: 3 :: R Regs Max Index (0 based
Pixel throughput (assuming 1 cycle texture lookup) 304.76 M

C:\Temp\NVShaderPerf_61_77>
```

# NVShaderPerf

**Direct3D Application**

HLSL ·····························▶

**DirectX Runtime**

Direct3D shader
op-codes ···················▶

**NVIDIA Driver**

DirectX Driver

API agnostic ·····▶
shader op-codes

Unified Compiler

················································▶ HW Binary

HW

# NVShaderPerf - Example

# FX Composer 1.5

- IDE for HLSL authoring, debugging and optimization
- Vertex and Pixel Shader scheduling
- Direct3D9 VS/PS op-code disassembly
- Advanced texture generation for baking Look Up Tables
- Visualization of RenderTargets



EverQuest® content courtesy Sony Online Entertainment

# FX Composer – 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 Perf**

| Untextured | p0 | Pixel Shader | GeForce 68( |

```
*******************************************
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

# Agenda

- Performance Tools Survey

- Practice

- Next generation Performance Tools

- Conclusion

- Q & A

# Basic Principles

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

# Pipelined Architecture (simplified view)

CPU → Geometry Storage → Geometry Processor → Rasterizer → Fragment Processor → Frame buffer

Texture Storage + Filtering → Fragment Processor

**Vertices**

**Pixels**

# The Terrible Bottleneck

Limits the speed of the pipeline

| CPU | → | Geometry Storage | → | Geometry Processor | → | Rasterizer | → | Fragment Processor | → | Frame buffer |

Texture Storage + Filtering

# Bottleneck Identification

- Need to identify it quickly and correctly
  - Guessing what it is without testing can waste a lot of coding time
- Two ways to identify a stage as the bottleneck
  - Modify the stage itself
  - Rule out the other stages

# Bottleneck Identification

- Modify the stage itself
  - By decreasing its workload



  - If performance 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



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

# Practice

- Now lets look at some sample problems and see if we can find out where the problem is
- Use NVPerfHUD to help

# Practice: Example 1

- A seemingly simple
  scene runs horribly slow
  - Narrow in on the bottleneck

# Practice: Example 1

- Dynamic vertex buffer
  - BAD creation flags

```
HRESULT hr = pd3dDevice->CreateVertexBuffer(
            6* sizeof( PARTICLE_VERT ),
            0,      //declares this as static
            PARTICLE_VERT::FVF,
            D3DPOOL_DEFAULT,
            &m_pVB,
            NULL );
```

# Practice: Example 1

- Dynamic vertex buffer
  - GOOD creation flags

```
HRESULT hr = pd3dDevice->CreateVertexBuffer(
            6* sizeof( PARTICLE_VERT ),
            D3DUSAGE_DYNAMIC |
            D3DUSAGE_WRITEONLY,
            PARTICLE_VERT::FVF,
            D3DPOOL_DEFAULT,
            &m_pVB,
            NULL );
```

# Practice: Example 1

- Dynamic Vertex Buffer
  - BAD Lock flags

m_pVB->Lock(0, 0,(void**)&quadTris, 0);

- No flags at all!?
  - That can't be good....

# Practice: Example 1

- Dynamic Vertex Buffer
  - GOOD Lock flags

```
m_pVB->Lock(0, 0,(void**)&quadTris,
D3DLOCK_NOSYSLOCK | D3DLOCK_DISCARD);
```

- Use D3DLOCK_DISCARD the first time you lock a vertex buffer each frame
  - And again when that buffer is full
  - Otherwise just use NOSYSLOCK

# Practice:  Example 2

- Another slow scene
  - What's the problem here

# Practice:  Example 2

- Texture bandwidth overkill
  - Use mipmaps
  - Use dxt1 if possible
    - Some cards can store compressed data in cache
  - Use smaller textures when they are fine
    - Does the grass blade really need a 1024x1024 texture?
      - Maybe

# Practice:  Example 3

- Another slow scene
  - Who wants a prize?

# Practice:  Example 3

- Expensive pixel shader
  - Can have huge performance effect
  - Only 3 verts, but maybe a million pixels
    - That's only 1024x1024

Look at all the pixels!!

# Practice: Example 3

32 cycles BAD ➡

**Shader Perf**  ✕

| TestFXCheapVSE ▾ | p0 ▾ | Pixel Shader ▾ | GeForce 6800 Ul ▾ |

```
********************************************
Target: GeForce 6800 Ultra (NV40) :: Unified Compiler: v61.77
Cycles: 31.50 :: R Regs Used: 7 :: R Regs Max Index (0 based): 6
Pixel throughput (assuming 1 cycle texture lookup) 206.45 MP/s
============================================
Shader performance using all FP16
Cycles: 20.00 :: R Regs Used: 3 :: R Regs Max Index (0 based): 2
Pixel throughput (assuming 1 cycle texture lookup) 320.00 MP/s
============================================
Shader performance using all FP32
Cycles: 31.50 :: R Regs Used: 7 :: R Regs Max Index (0 based): 6
Pixel throughput (assuming 1 cycle texture lookup) 206.45 MP/s
********************************************
PS Instructions: 35
ps_2_0
def c0, 0, 1, 0, 0
dcl t0.xy
dcl t1.xyz
dcl t2.xyz
dcl t3
dcl t4
dcl t5
dcl_2d s0
dp4 r0.w, t4, t4
rsq r0.w, r0.w
mul r0.xyz, r0.w, t4
mul r1 xyz r0 u t2
```

✉ Properties    📂 Shader Perf

# Practice: Example 3

12 cycles GOOD

**Shader Perf**

TestFXCheapVSM ▾ | p0 ▾ | Pixel Shader ▾ | GeForce 6800 Ul ▾

```
*******************************************
Target: GeForce 6800 Ultra (NV40) :: Unified Compiler: v61.77
Cycles: 12.00 :: R Regs Used: 2 :: R Regs Max Index (0 based): 1
Pixel throughput (assuming 1 cycle texture lookup) 533.33 MP/s
===========================================
Shader performance using all FP16
Cycles: 12.00 :: R Regs Used: 2 :: R Regs Max Index (0 based): 1
Pixel throughput (assuming 1 cycle texture lookup) 533.33 MP/s
===========================================
Shader performance using all FP32
Cycles: 11.00 :: R Regs Used: 3 :: R Regs Max Index (0 based): 2
Pixel throughput (assuming 1 cycle texture lookup) 581.82 MP/s
*******************************************
PS Instructions: 20
ps_2_0
def c0, 0, 1, 0, 0
dcl t0.xy
dcl t3
dcl t6
dcl t7
dcl_2d s0
dp4 r0.w, t3, t3
rsq r0.w, r0.w
mul_pp r1, r0.w, t3
dp4_pp r4.w, r1, t6
add_pp r0.w, r4.w, r4.w
mad_pp r0, r0.w, r1, t6
```
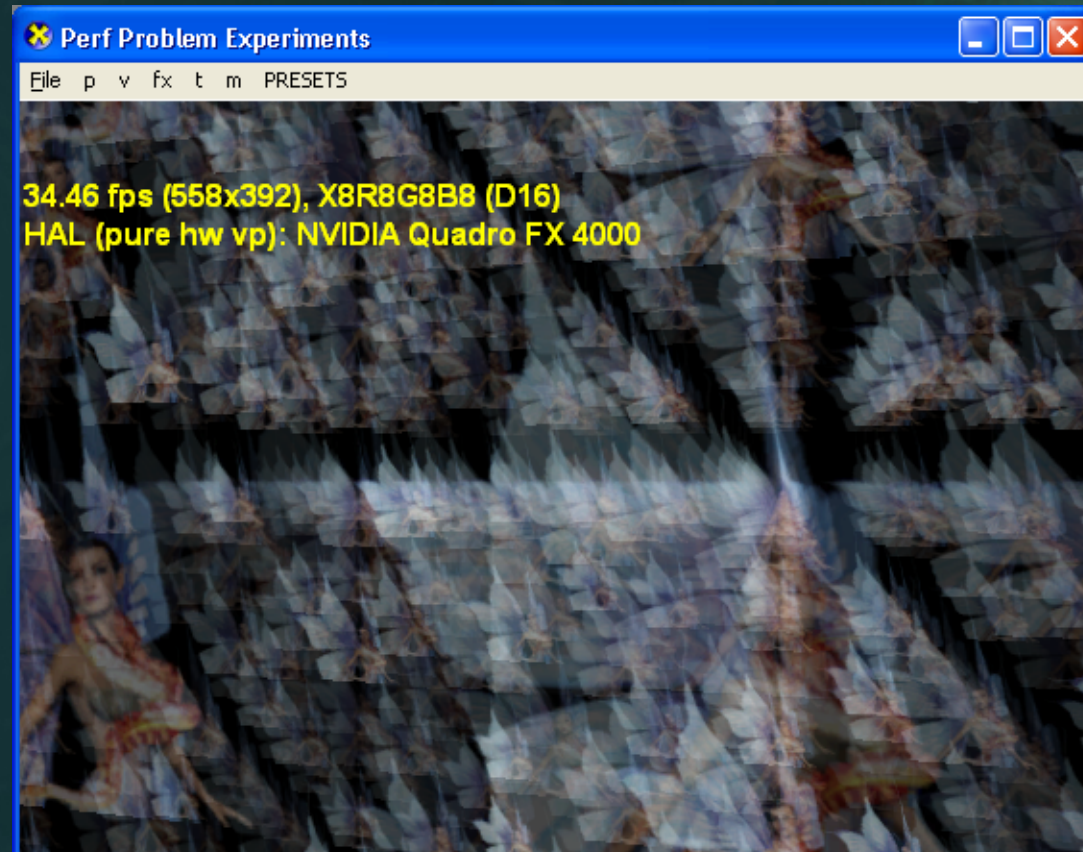
▢ Properties  ▢ Shader Perf

# Practice: Example 3

- ## What changed?

  - ### Moved math that was constant across the triangle into the vertex shader

  - ### Used 'half' instead of 'float'

  - ### Got rid of normalize where it wasn't necessary

    - See Normalization Heuristics
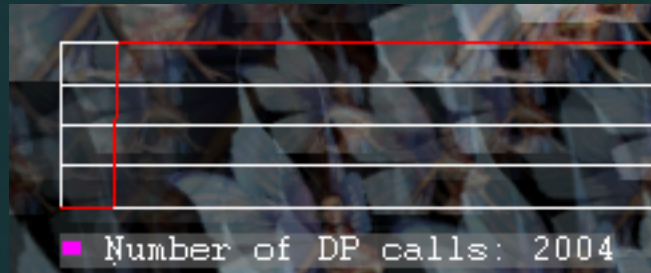    - http://developer.nvidia.com

# Practice:  Example 4

- The last one
  - Audience: there are no more prizes, but we've locked the doors

# Practice:  Example 4

- Too many batches
  - Was sending every quad as it's own batch
  - Instead, group quads into one big VB then send that with one call



Number of DP calls: 2004

# Practice: Example 4

- What if they use different textures?
  - Use texture atlases
  - Put the two textures into a single texture and use a vertex and pixel shader to offset the texture coordinates

# Agenda

- Performance Tools Survey
- Performance Methodologies and Practice
- **Next generation Performance Tools**
- Conclusion
- Q & A

# Next Generation Performance Tools

- NVIDIA Performance Kit (PerfKit)
  - Instrumented Driver
  - NVIDIA Developer Control Panel (NVDevCPL)
  - NVIDIA Plug-in for Microsoft PIX for Windows
  - Sample Code for DirectX
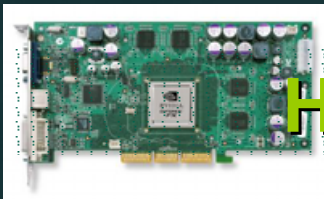
# Problem

**Application**

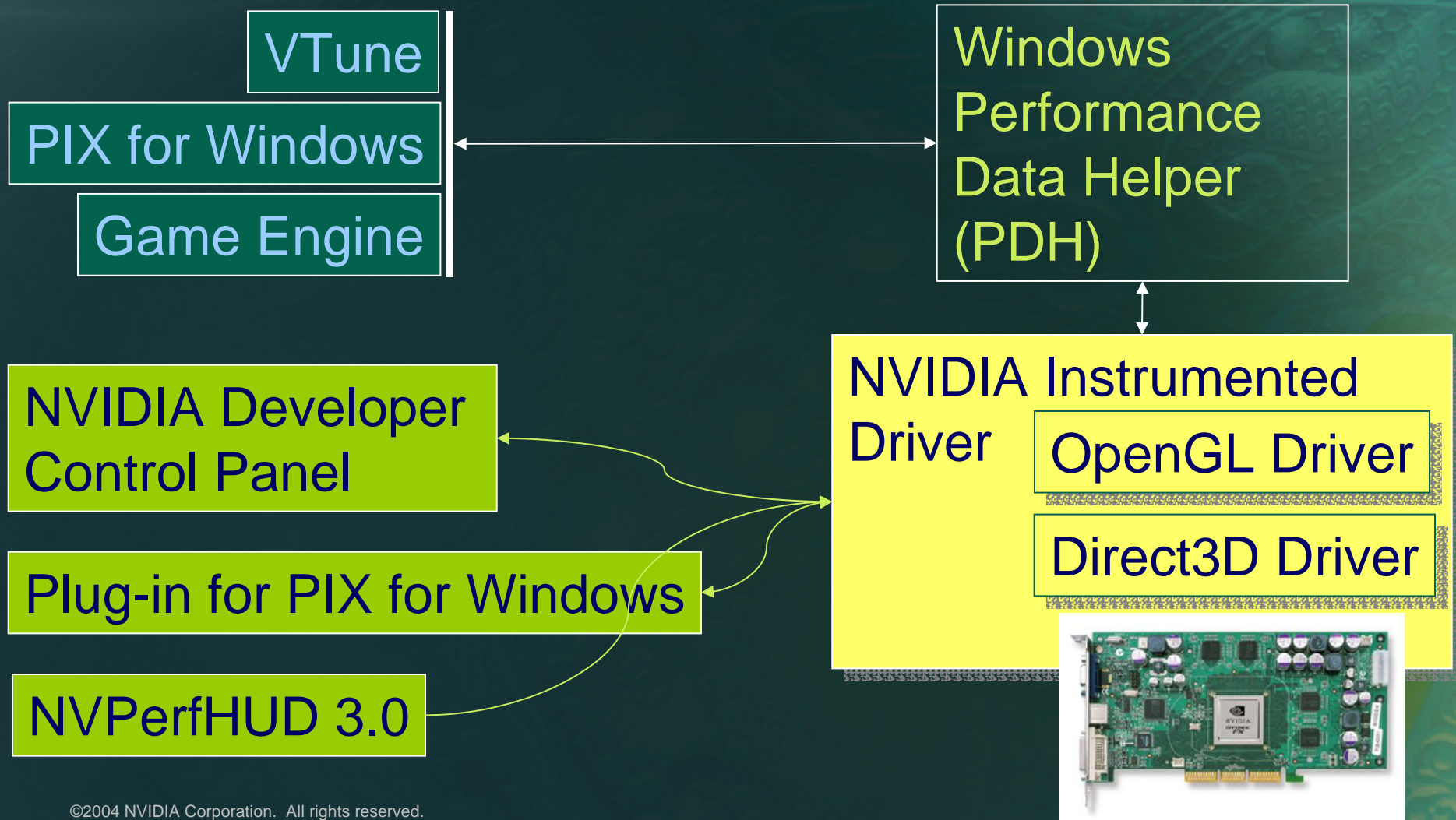**Common Profilers**

*Microsoft* **DIRECTX** API

**PIX for Windows**

**Driver**

**Hardware**

**How to evaluate performance here?**

# Solution

VTune

PIX for Windows

Game Engine

Windows Performance Data Helper (PDH)

NVIDIA Developer Control Panel

Plug-in for PIX for Windows

NVPerfHUD 3.0

NVIDIA Instrumented Driver
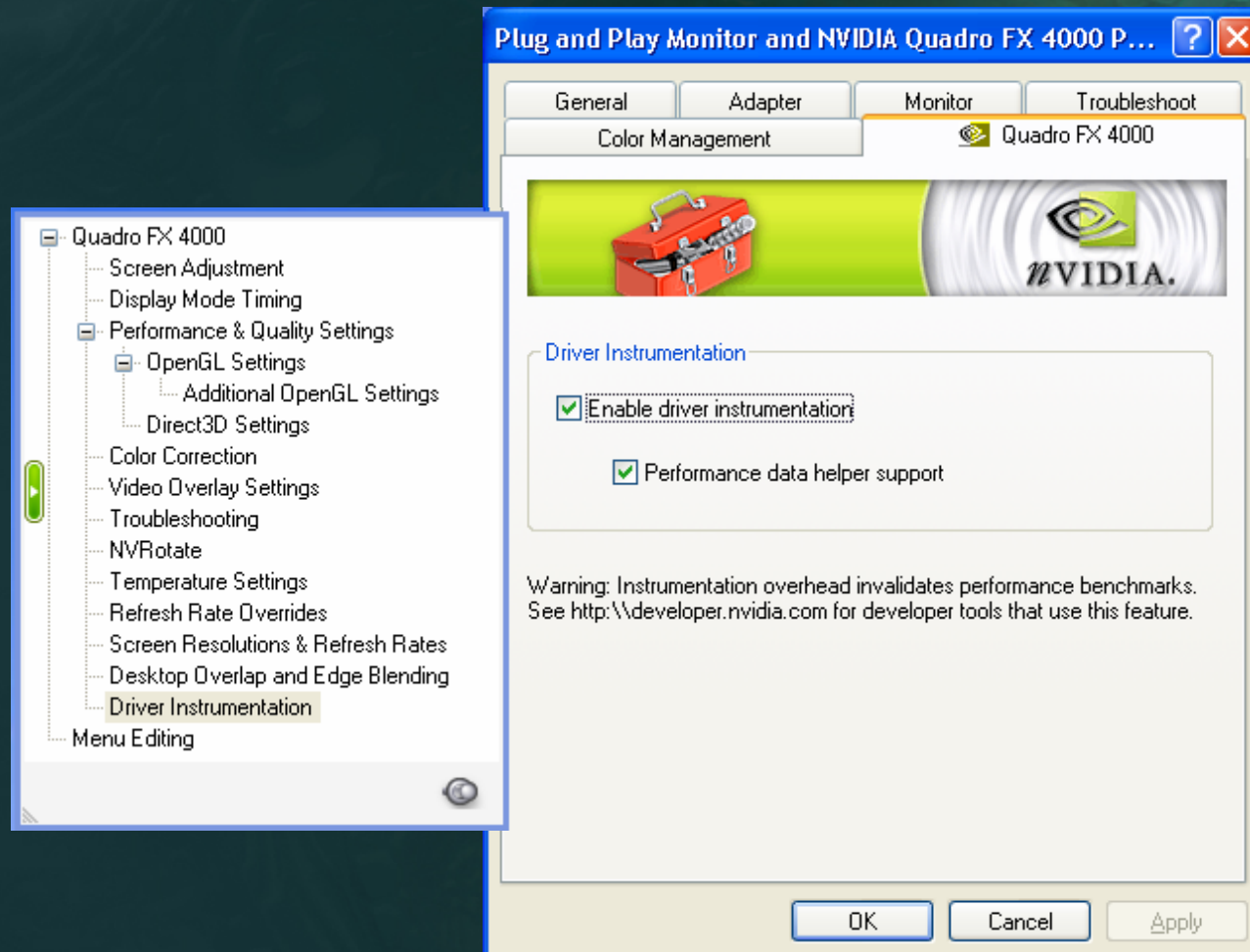
OpenGL Driver

Direct3D Driver

# Instrumented Driver

- Special Instrumented Driver
    - Built with regular drivers
    - Includes NVPMAPI.DLL
- Exposes Driver and HW Performance Counters
- Compatible with Windows WMI and PDH
- New Driver Instrumentation tab in NVIDIA Display Control Panel
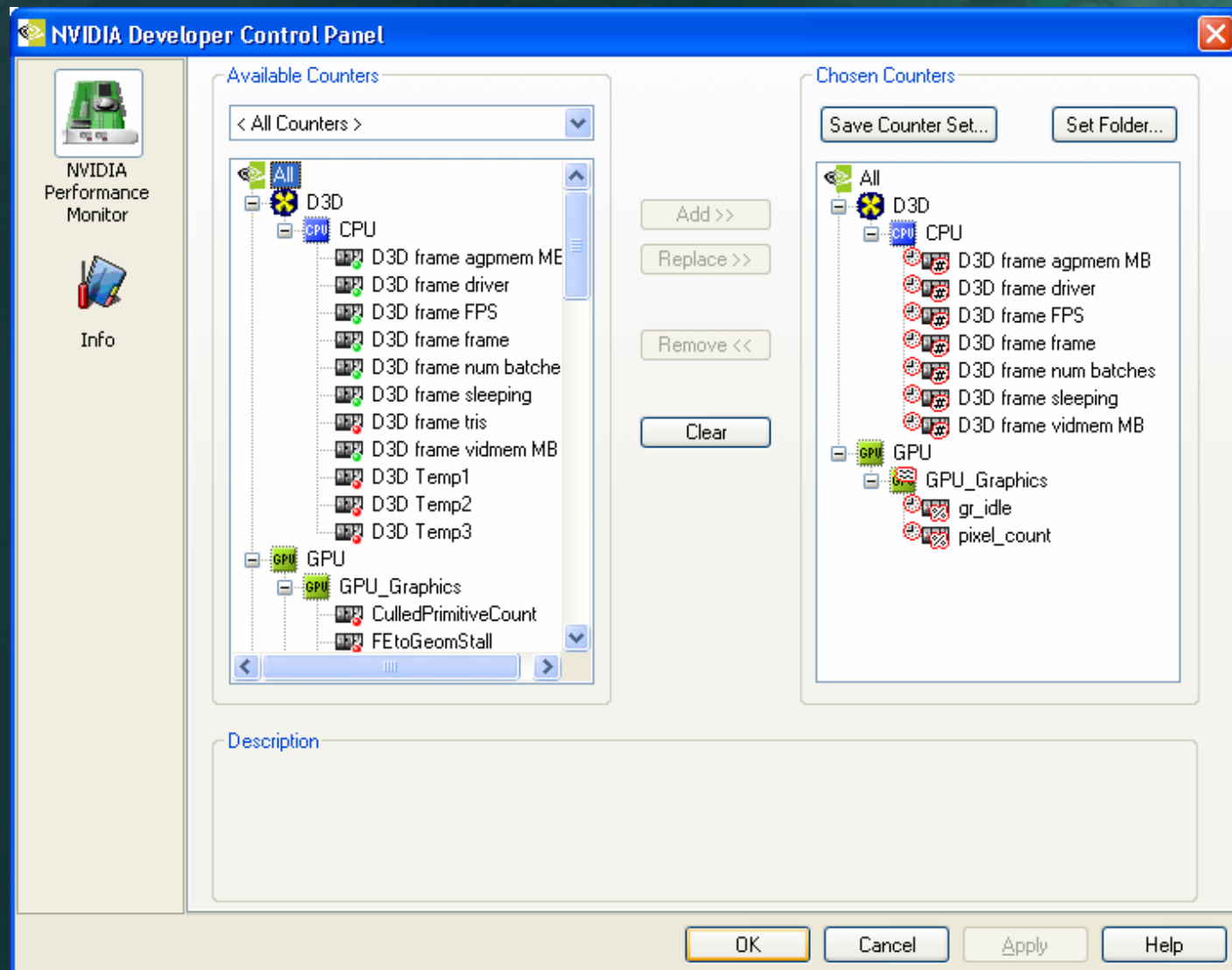
# Instrumented Driver

# Developer Control Panel

- Control per-counter specifics
  - Enabled or not
  - Raw values or % values
  - Etc.
- Manage multiple counter-sets
  - Tray Icon: fast application of presets
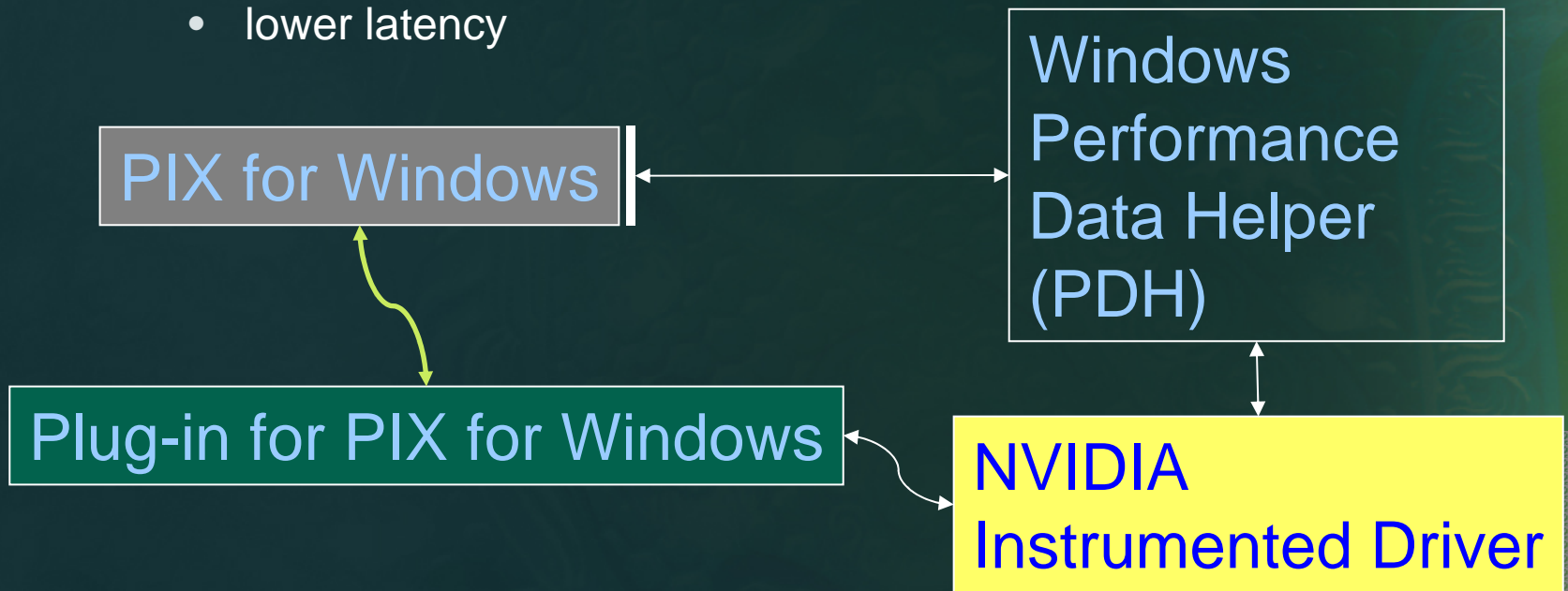- Provides HW specific information
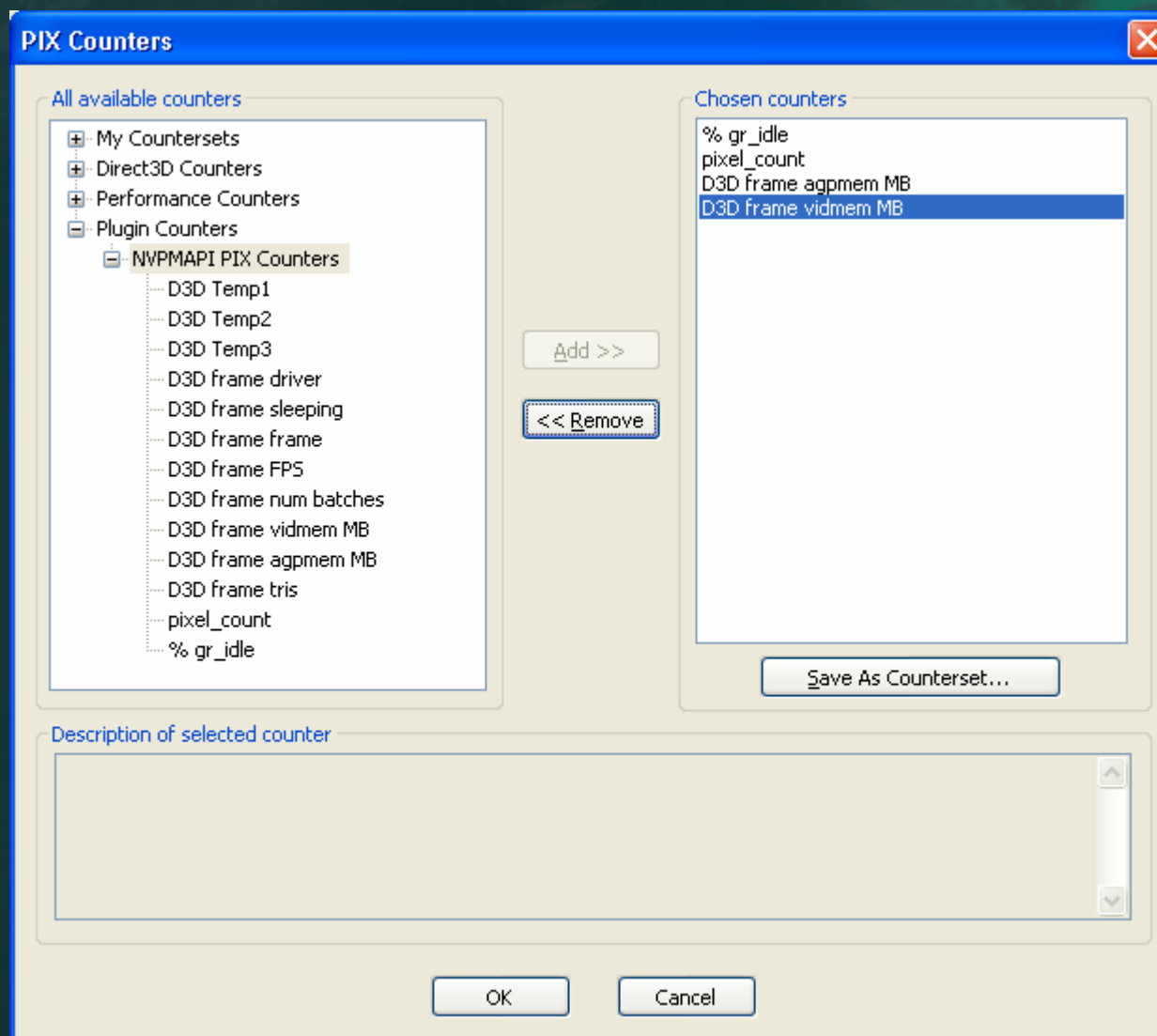
# Developer Control Panel: Demo

# PIX for Windows Plug-in

- NVIDIA's performance counters
  - PIX's PDH adaptor
  - NVIDIA's Pix Plug-in
    - higher frequency
    - lower latency

PIX for Windows

Plug-in for PIX for Windows

Windows Performance Data Helper (PDH)

NVIDIA Instrumented Driver

# Plug-in for PIX: Demo

# NVPerfKit Code Samples

- Includes C++ helper classes for PDH access and display
  - PDHHelper
  - Trace<T>
  - TraceDisplay
    - Various display types
    - Direct3D implementation
- Sample Code and App
  - Illustrates sampling issues and dynamic reconfigurability

# Conclusion

- Comprehensive Suite of Performance Tools
  - performance information at all levels
    - Direct3D API
    - Direct3D Driver
    - Hardware
- Provided in a variety of venues
  - Microsoft WMI/PDH
  - Microsoft PIX for Windows
  - User application
  - NVPerfHUD

# Questions?

- What else can we do for you?
  - [sdkfeedback@nvidia.com](mailto:sdkfeedback@nvidia.com)