



DirectX 10 Performance

Ignacio Llamas, NVIDIA

illamas@nvidia.com

Nicolas Thibieroz, AMD

nicolas.thibieroz@amd.com

WWW.GDCONF.COM



Outline

- ④ General DX10 API usage
 - ④ Designed for performance
 - ④ Batching and Instancing
 - ④ State Management
 - ④ Constant Buffer Management
 - ④ Resource Updates and Management
 - ④ Reading the Depth Buffer
 - ④ MSAA
- ④ Optimizing your DX10 Game
 - ④ or how to work around GPU bottlenecks
- ④ IHV-specific optimizations



Game Developers
Conference 08

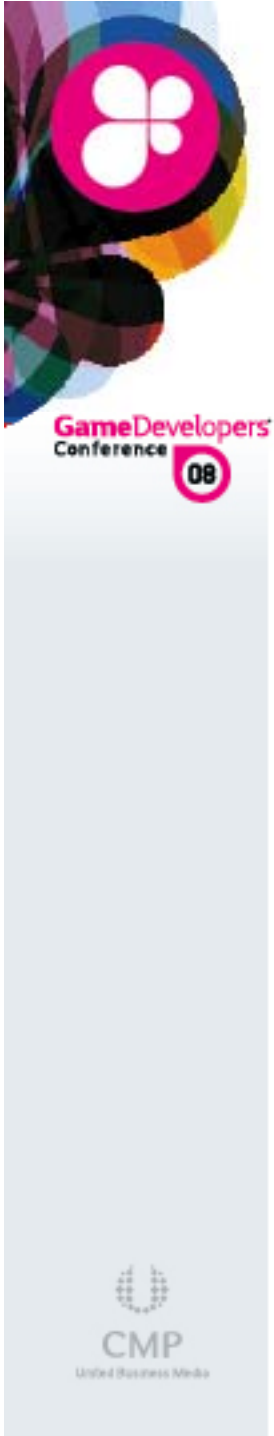
Color Guide For IHV-specific Advice

 AMD

 NVIDIA


CMP
United Business Media

WWW.GDCONF.COM



DX10 Runtime and Driver. Designed for Performance

- ⊕ DX10 validation moved from runtime to creation time
 - ⊕ Only basic error checking at runtime
- ⊕ Immutable state objects
 - ⊕ Can be pre-computed and cached
 - ⊕ Subset of command buffer at creation time
- ⊕ Vista driver model delegates scheduling and memory management to OS
 - ⊕ Pro: more responsive system, GPU sharing across apps
 - ⊕ Con: harder to guarantee performance if multiple apps share the GPU
 - ⊕ Fullscreen mode should be fine



Game Developers
Conference 08

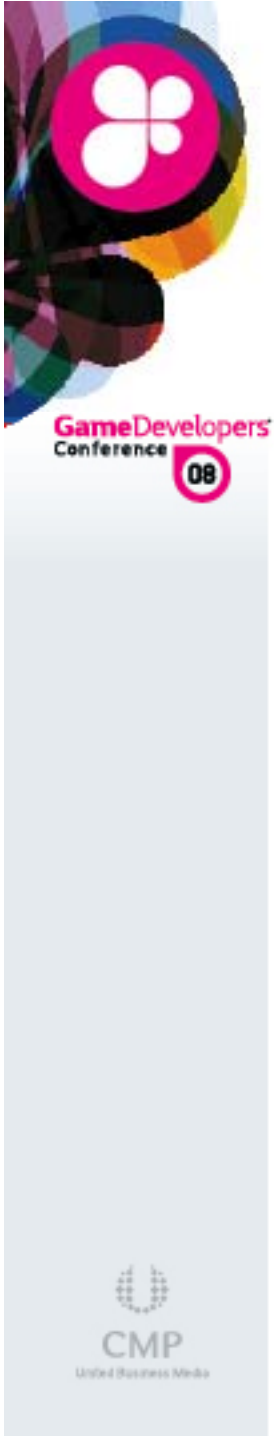
Batch Performance

- ④ The truth about DX10 batch performance
- ④ “Simple” porting job will not yield expected performance
- ④ Need to use DX10 features to yield gains:
 - ④ Geometry instancing or batching
 - ④ Intelligent usage of state objects
 - ④ Intelligent usage of constant buffers
 - ④ Texture arrays



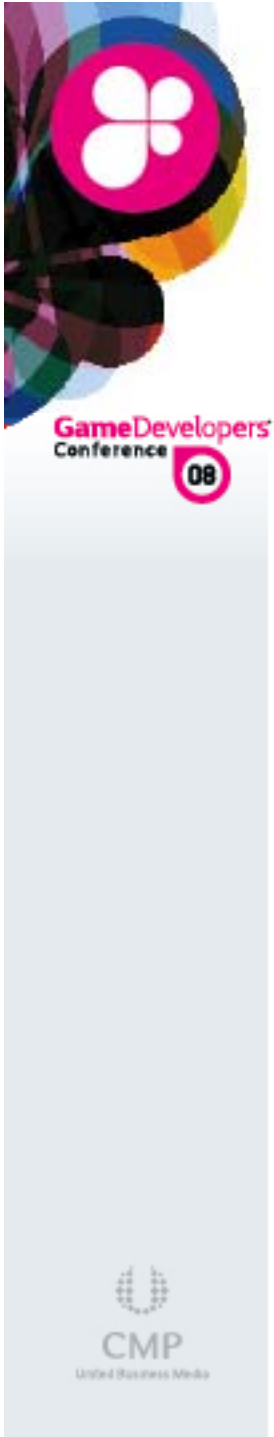
Geometry Instancing

- ⊕ Better instancing support in DX10
- ⊕ Use "System Values" to vary rendering
 - ⊕ `SV_InstanceID`, `SV_PrimitiveID`, `SV_VertexID`
 - ⊕ Additional streams not required
 - ⊕ Pass these to PS for texture array indexing
 - ⊕ Highly-varied visual results in a single draw call
- ⊕ Watch out for:
 - ⊕ Texture cache trashing if sampling textures from system values (`SV_PrimitiveID`)
 - ⊕ Too many attributes passed from VS to PS
 - ⊕ **InputAssembly** bottlenecks due to instancing
 - ⊕ **Solution: Load() per-instance data from Buffer in VS or PS using `SV_InstanceID`**



State Management

- ④ DX10 uses immutable “state objects”
 - ④ Input Layout Object
 - ④ Rasterizer Object
 - ④ DepthStencil Object
 - ④ Blend Object
 - ④ Sampler Object
- ④ DX10 requires a new way to manage states
 - ④ A naive DX9 to DX10 port *will* cause problems here
 - ④ **Always create state objects at load-time**
 - ④ Avoid duplicating state objects
 - ④ Recommendation to sort by states still valid in DX10!

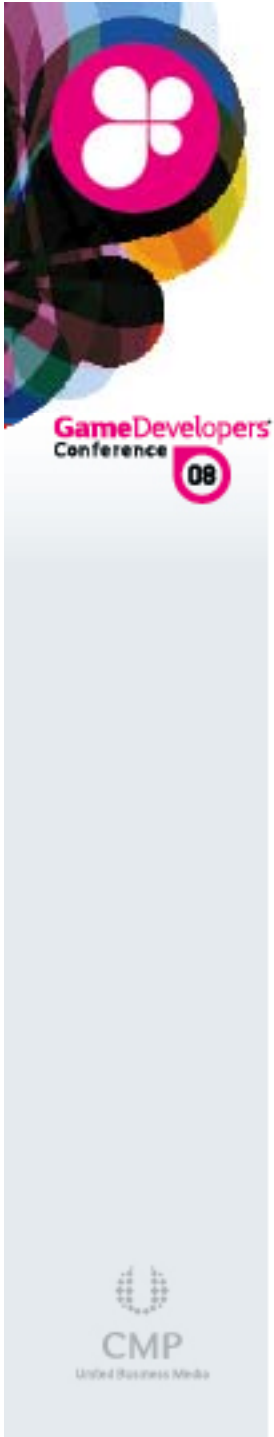


Constant Buffer Management (1)

- ⊗ Probably a major cause of poor performance in initial naïve DX10 ports!
- ⊗ Constants are declared in buffers in DX10

```
cbuffer PerFrameConstants  cbuffer SkinningMatricesConstants
{
    float4x4  mView;
    float      fTime;
    float3     fWindForce;
    // etc.
};
```

- ⊗ When any constant in a cbuffer is updated the full cbuffer has to be uploaded to GPU
- ⊗ Need to strike a good balance between:
 - ⊗ Amount of constant data to upload
 - ⊗ Number calls required to do it (== # of cbuffers)



Constant Buffer Management (2)

- ④ Use a pool of constant buffers *sorted by frequency of updates*
- ④ Don't go overboard with number of cbuffers!
 - ④ (3-5 is good)
- ④ Sharing cbuffers between shader stages can be a good thing
- ④ Example cbuffers:
 - ④ PerFrameGlobal (time, per-light properties)
 - ④ PerView (main camera xforms, shadowmap xforms)
 - ④ PerObjectStatic (world matrix, static light indices)
 - ④ PerObjectDynamic (skinning matrices, dynamic lightIDs)



Constant Buffer Management (3)

- ③ Group constants by access pattern to help cache reuse due to locality of access
- ③ Example:

```
float4 PS_main(PSInput in)
{
    float4 diffuse = tex2D0.Sample(mipmapSampler, in.Tex0);
    float ndotl = dot(in.Normal, vLightVector.xyz);
    return ndotl * vLightColor * diffuse;
}
```

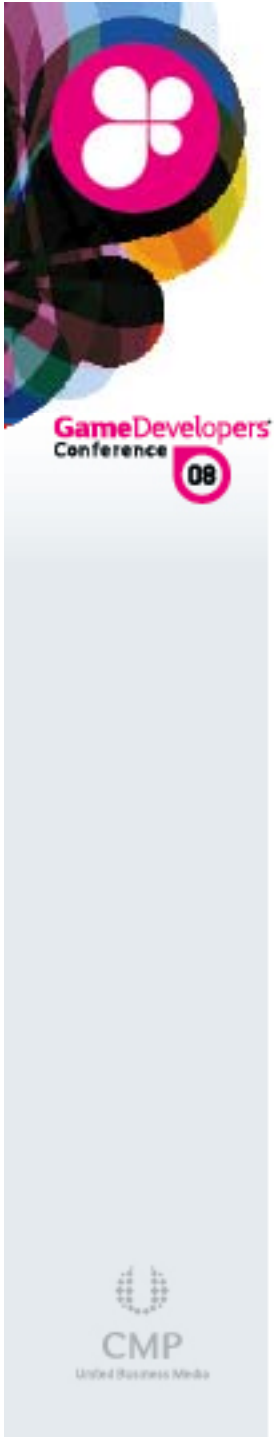
```
cbuffer PerFrameConstants
{
    float4    vLightVector;
    float4    vLightColor;
    float4    vOtherStuff[32];
};
GOOD
```

```
cbuffer PerFrameConstants
{
    float4    vLightVector;
    float4    vOtherStuff[32];
    float4    vLightColor;
};
BAD
```



Constant Buffer Management (4)

- ④ Careless DX9 port results in a single `$Globals cbuffer` containing all constants, many of them unused
- ④ `$Globals cbuffer` typically yields bad performance:
 - ④ Wasted CPU cycles updating unused constants
 - ④ Check if used: `D3D10_SHADER_VARIABLE_DESC.uFlags`
 - ④ cbuffer contention
 - ④ Poor cbuffer cache reuse due to suboptimal layout
- ④ When compiling SM3 shaders for SM4+ target with `D3D10_SHADER_ENABLE_BACKWARDS_COMPATIBILITY`: use conditional compilation to declare `cbuffers` (e.g. `#ifdef DX10 cbuffer{ #endif`)



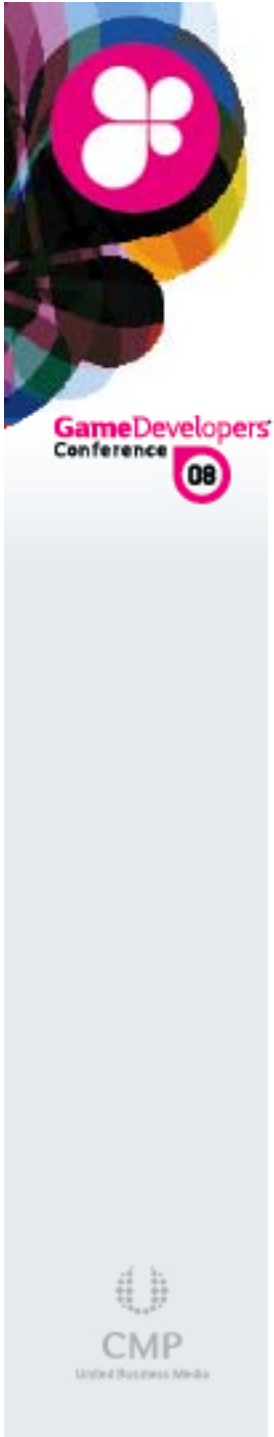
Constant Buffer Management (5)

- ④ Consider **tbuffer** if access pattern is more random than sequential
 - ④ **tbuffer** access uses texture Loads, so higher latency but higher performance sometimes
 - ④ Watch out for texture-bound cases resulting from **tbuffer** usage
- ④ Use **tbuffer** if you need more data in a single buffer
 - ④ **cbuffer** limited to 4096*128-bit
 - ④ **tbuffer** limited to 128 megabytes



Resource Updates

- ⊗ In-game destruction and creation of Texture and Buffer resources has a significant impact on performance:
 - ⊗ Memory allocation, validation, driver checks
- ⊗ Create all resources up-front if possible
 - ⊗ During level load, cutscenes, or any non-performance critical situations
- ⊗ At runtime: replace contents of existing resources, rather than destroying/creating new ones



Resource Updates: Textures

- ④ Avoid **UpdateSubresource()** for textures
 - ④ Slow path in DX10
(think **DrawPrimitiveUP()** in DX9)
 - ④ Especially bad with larger textures!
- ④ Use ring buffer of intermediate **D3D10_USAGE_STAGING** textures
 - ④ Call **Map(D3D10_MAP_WRITE, ...)** with **D3D10_MAP_FLAG_DO_NOT_WAIT** to avoid stalls
 - ④ If Map fails in all buffers: either stall waiting for Map or allocate another resource (cache warmup time)
 - ④ Copy to textures in video memory (**D3D10_USAGE_DEFAULT**):
 - ④ **CopyResource()** or **CopySubresourceRegion()**



Resource Updates: Buffers

- ④ To update a Constant buffer
 - ④ `Map(D3D10_MAP_WRITE_DISCARD, ...)` ;
 - ④ `UpdateSubResource()`
 - ④ Recall full buffer must be updated, but with `Map()` CPU can skip parts that the shader does not care about. All the data must be uploaded to GPU though

- ④ To update a dynamic Vertex/Index buffer
 - ④ Use a *large* shared ring-buffer type; writing to unused portions of buffer using:
 - ④ `Map(D3D10_MAP_WRITE_DISCARD, ...)` when full or if possible the first time it is mapped at every frame
 - ④ `Map(D3D10_MAP_WRITE_NO_OVERWRITE, ...)` thereafter
 - ④ Avoid `UpdateSubResource()`
 - ④ not as good as `Map()` in this case either



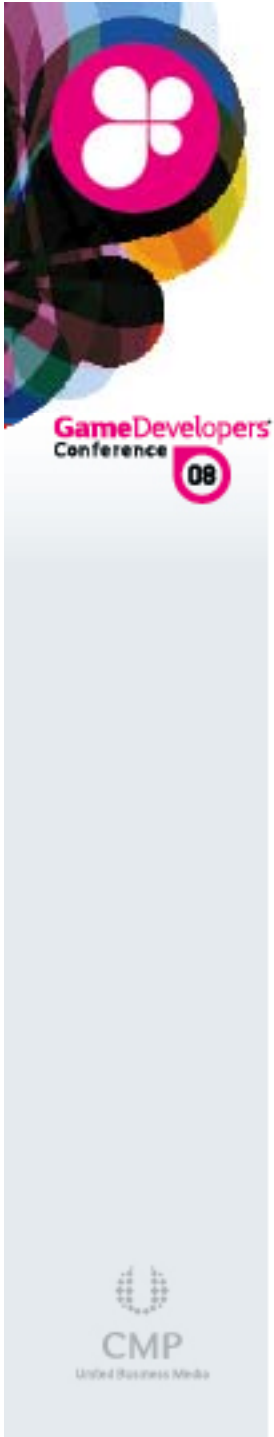
Game Developers
Conference

08

Accessing Depth and Stencil

- ④ DX10 enables the depth buffer to be read back as a texture
- ④ Enables features without requiring a separate depth render
 - Atmosphere pass
 - Soft particles
 - Depth of Field
 - Deferred shadow mapping
 - Screen-space ambient occlusion
 - Etc.
- ④ Popular features in most recent game engines





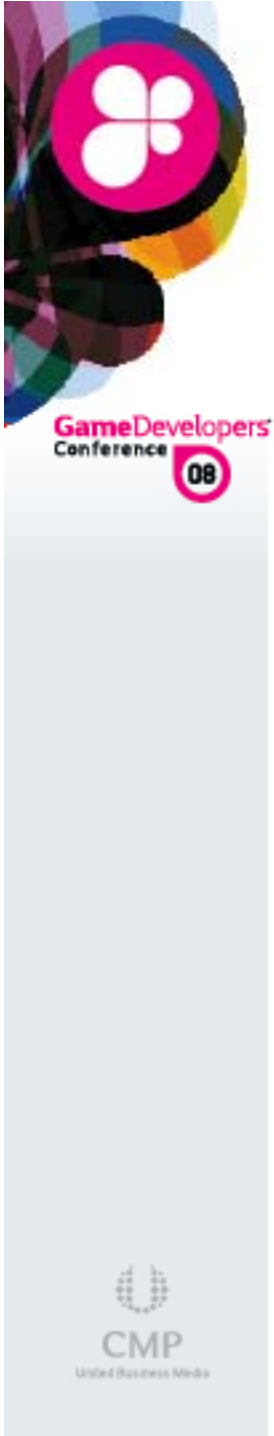
Accessing Depth and Stencil with MSAA

- ④ DX10.0: reading a depth buffer as SRV is only supported in single sample mode
 - ④ Requires a separate render path for MSAA
- ④ Workarounds:
 - ④ Store depth in alpha of main FP16 RT
 - ④ Render depth into texture in a depth pre-pass
 - ④ Use a secondary rendertarget in main color pass
- ④ DX10.1 allows depth buffer access as Shader Resource View in all cases:
 - ④ Fewer shaders
 - ④ Smaller memory footprint
 - ④ Better orthogonality



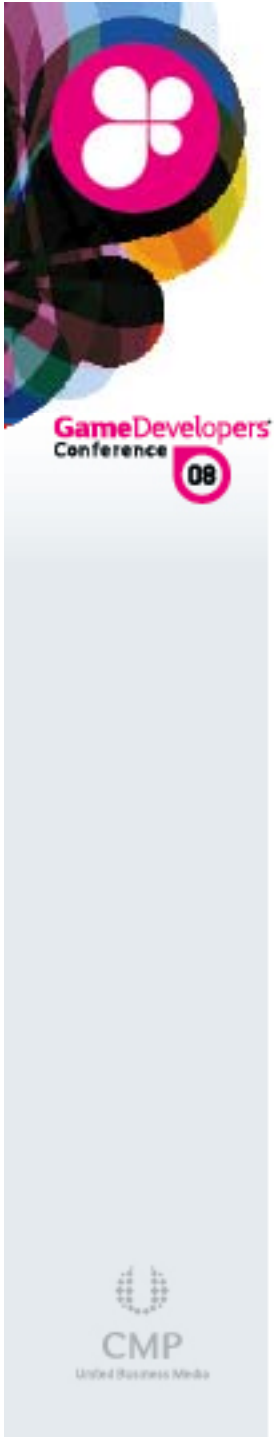
MultiSampling Anti-Aliasing

- ⊗ MSAA resolves cost performance
 - ⊗ Cost varies across GPUs but it is never free
 - ⊗ Avoid redundant resolves as much as possible
 - E.g.: no need to perform most post-process ops on MSAA RT. Resolve once, *then* apply p.p. effects
- ⊗ No need to allocate SwapChain as MSAA
 - ⊗ Apply MSAA only to rendertargets that matter
- ⊗ Be aware of CSAA on NVIDIA hardware:
 - Certain `DXGI_SAMPLE_DESC.Quality` values will enable higher-quality but slightly costlier MSAA mode
 - See <http://developer.nvidia.com/object/coverage-sampled-aa.html>



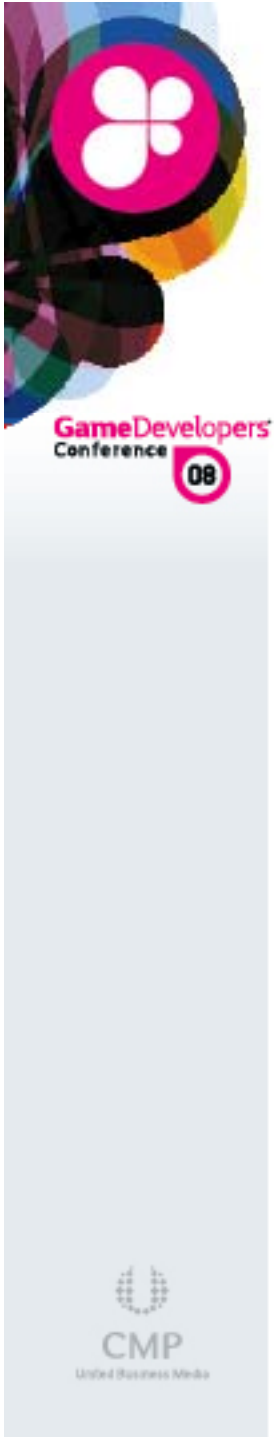
Optimizing your DX10 Game

- ④ Use **PerfHUD** or **GPUPerfStudio** to identify bottlenecks:
 - ④ Step 1: are you GPU or CPU bound
 - ④ Check GPU idle time
 - ④ If GPU is idle you are probably CPU bound either by other CPU workload on your application or by CPU-GPU synchronization
 - ④ Step 2: if GPU bound, identify the top buckets and their bottlenecks
 - ④ Use PIX or **PerfHUD** Frame Profiler for this
 - ④ Step 3: try to reduce the top bottleneck/s



If Input Assembly is the bottleneck

- ④ Optimize IB and VB for cache reuse
 - ④ Use `ID3DXMesh::Optimize()` or other tools
- ④ Reduce number of vector attributes
 - ④ Pack several scalars into single 4-scalar vector
- ④ Reduce vertex size using packing tricks:
 - ④ Pack normals into a float2 or even RGBA8
 - ④ Calculate binormal in VS
 - ④ Use lower-precision formats
- ④ Use reduced set of VB streams in shadow and depth-only passes
 - ④ Separate position and 1 texcoord into a stream
 - ④ Improves cache reuse in pre-transform cache
 - ④ Also use shortest possible shaders



If Vertex Shader is the bottleneck

- ④ Improve culling and LOD (also helps IA):
 - ④ Look at wireframe in debugging tool and see if it's reasonable
 - ④ Check for percentage of triangles culled:
 - ④ Frustum culling
 - ④ Zero area on screen
 - ④ Use other scene culling algorithms
 - ④ CPU-based culling
 - ④ Occlusion culling
- ④ Use Stream-Output to cache vertex shader results for multiple uses
 - ④ E.g.: StreamOut skinning results, then render to shadowmap, depth prepass and shading pass
 - ④ StreamOut pass writes point primitives (vertices) Same index buffer used in subsequent passes



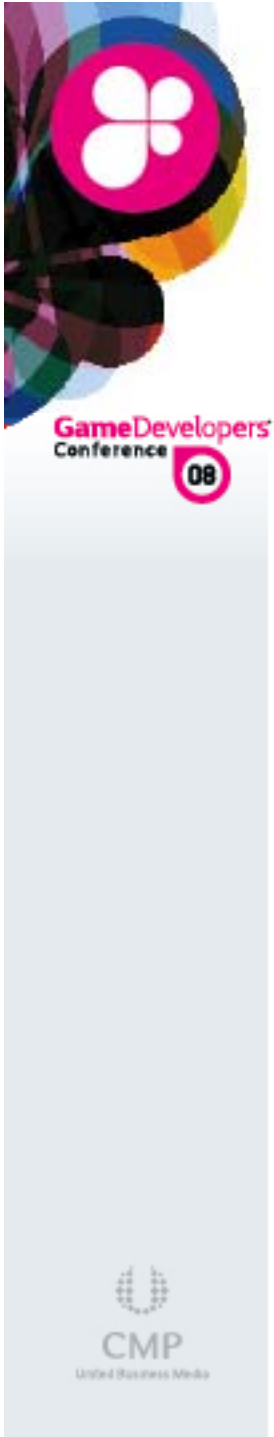
If Geometry Shader is the bottleneck

- ⊗ Make sure *maxvertexcount* is as low as possible
 - ⊗ *maxvertexcount* is a shader constant declaration → need different shaders for different values
 - ⊗ Performance drops as output size increases
- ⊗ Minimize the size of your output and input vertex structures
- ⊗ GS not designed for large-expansion algorithms like tessellation
 - ⊗ Due to required ordering and serial execution
 - ⊗ See Andrei Tatarinov's talk on Instanced Tessellation
- ⊗ Consider using instancing in current hardware
- ⊗ Move some computation to VS to avoid redundancy
- ⊗ **Keep GS shaders short**
- ⊗ **Free ALUs in GS because of latency**
 - ⊗ Can be used to cull geometry (backface, frustum)



If Stream-Output is the bottleneck

- ⌚ Avoid reordering semantics in the output declaration
 - ⌚ Keep them in same order as in output structure
- ⌚ You may have hit bandwidth limit
 - ⌚ SO bandwidth varies by GPU
- ⌚ Remember you don't need to use a GS if you are just processing vertices
 - ⌚ Use ConstructGSWithSO on Vertex Shader
- ⌚ Rasterization can be used at the same time
 - ⌚ Only enable it if needed (binding RenderTarget)



If Pixel Shader is the bottleneck (1)

- ④ Verify by replacing with simplest PS
 - ④ Use **PerfHUD** / **GPUPerfStudio**
- ④ Move computations to Vertex Shader
- ④ Use pixel shader LOD
- ④ Only use **discard** or **clip()** when required
- ④ **discard** or **clip()** as early as possible
 - ④ GPU can skip remaining instructions if test succeeds
- ④ Use common app-side solutions to maximize pixel culling efficiency:
 - ④ Depth prepass (most common)
 - ④ Render objects front to back
 - ④ Triangle sort to optimize both for post-transform cache and Z culling within a single mesh
 - ④ Stencil/scissor/user clip planes to tag shading areas
 - ④ Deferred shading



If Pixel Shader is the bottleneck (2)

- ③ Shading can be avoided by Z/Stencil culling
 - ③ Coarse (ZCULL / Hi-Z)
 - ③ Fine-grained (EarlyZ)

- ③ Coarse Z culling is transparent, but it **may underperform** if:
 - ③ If shader writes depth
 - ③ High-frequency information in depth buffer
 - ③ If you don't clear the depth buffer using a "clear" (avoid clearing using fullscreen quads)



GameDevelopers
Conference

08

If Pixel Shader is the bottleneck (3)

- ⌚ Fine-grained Z culling is not always active

- ⌚ Disabled on current hardware if:

- ⌚ PS writes depth (SV_Depth)

- ⌚ Z or Stencil writes combined with:

- ⌚ Alpha test is enabled (DX9 only)

- ⌚ discard / texkill in shaders

- ⌚ AlphaToCoverageEnable = true

- ⌚ Disabled on current NVIDIA HW if:

- ⌚ PS reads depth (.z) from SV_Position input

- ⌚ Use .w (view-space depth) if possible

- ⌚ Z or Stencil writes combined with:

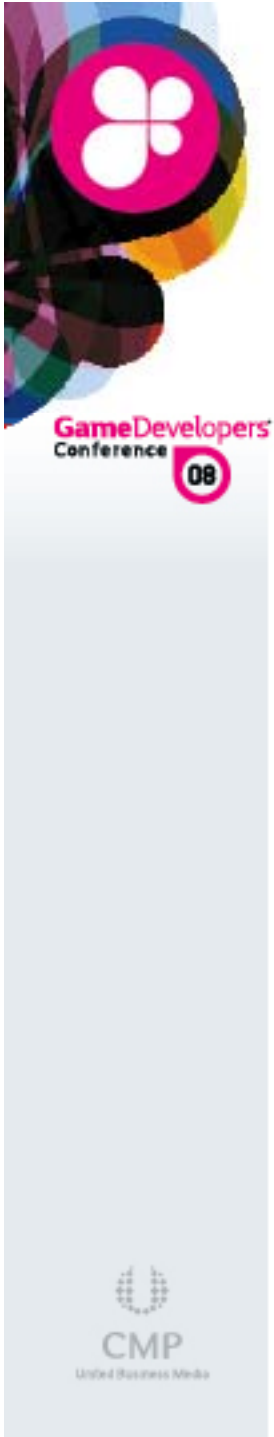
- ⌚ Samplemask != 0xffffffff



CMP

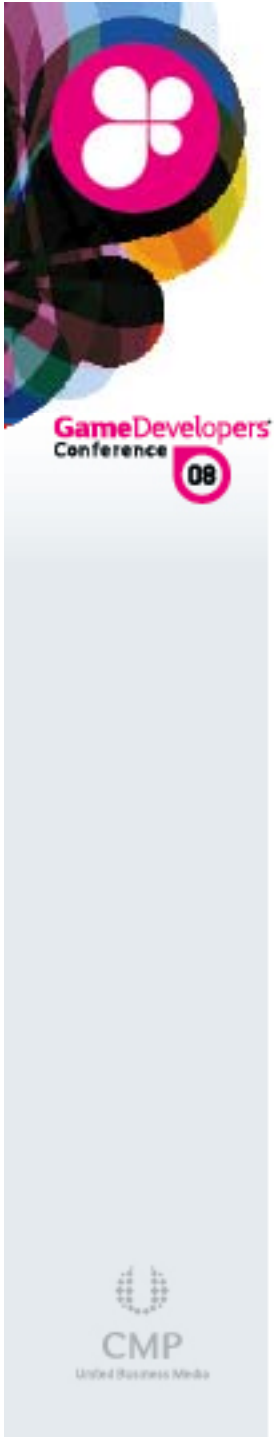
United Business Media

WWW.GDCONF.COM



Any Shader is *still* the bottleneck (1)

- ④ Use IHV tools:
 - ④ AMD: GPUShaderAnalyzer
 - ④ NVIDIA: ShaderPerf
- ④ Be aware of appropriate ALU to TEX **hardware** instruction ratios:
 - ④ 4 5D-vector ALU per TEX on AMD
 - ④ 10 scalar ALU per TEX on NVIDIA GeForce 8 series
- ④ Check for excessive register usage
 - ④ > 10 vector registers is high on GeForce 8 series
 - ④ Simplify shader, disable loop unrolling
 - ④ DX compiler behavior may unroll loops so check output
- ④ Use dynamic branching to skip instructions
 - ④ Make sure branching has high coherency



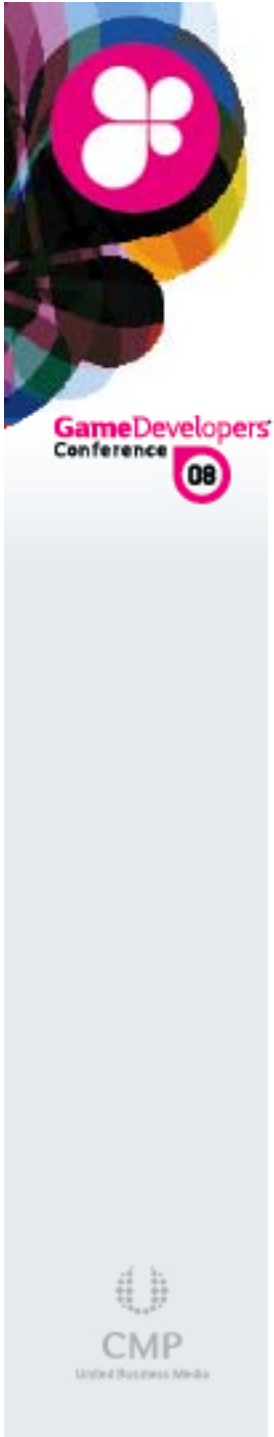
Any Shader is *still* the bottleneck (2)

- ⊗ Some instructions operate at a slower rate
 - ⊗ Integer multiplication and division
 - ⊗ Type conversion (float to int, int to float)
- ⊗ Too many of those can cause a bottleneck in your code
- ⊗ In particular watch out for type conversions
 - ⊗ Remember to declare constants in the same format as the other operands they're used with!



If Texture is the bottleneck (1)

- ④ Verify by replacing textures with 1x1 texture
 - ④ PerfHUD or GPUPerfStudio can do this
- ④ Basic advice:
 - ④ Enable mipmapping
 - ④ Use compressed textures where possible
 - ④ Block-compressed formats
 - ④ Compressed float formats for HDR
 - ④ Avoid negative LOD bias (aliasing != sharper)
- ④ If multiple texture lookups are done in a loop
 - ④ Unrolling partially may improve batching of texture lookups, reducing overall latency
 - ④ However this may increase register pressure
 - ④ Find the right balance



If Texture is the bottleneck (2)

- ③ DirectX compiler moves texture instructions that compute LOD out of branches
 - ③ Use SampleLevel (no anisotropic filtering)
 - ③ SampleGrad can be used too, but beware of the extra performance cost
- ③ Texture cache misses may be high due to poor coherence
 - ③ In particular in post-processing effects
 - ③ Modify access pattern
- ③ Not all textures are equal in sample performance
 - ③ Filtering mode
 - ③ Volume textures
 - ③ Fat formats (128 bits)
 - ③ **64-bit *integer* textures**



Game Developers
Conference

08

If ROP is the bottleneck: Causes

- ⊕ Pixel shader is too cheap ☺
- ⊕ Large pixel formats
- ⊕ High resolution
- ⊕ Blending
- ⊕ MSAA
- ⊕ MRT
- ⊕ Rendering to system memory over PCIe (parts with no video memory)
- ⊕ Typical problem with particle effects: little geometry, cheap shading, but high overdraw using blending



CMP

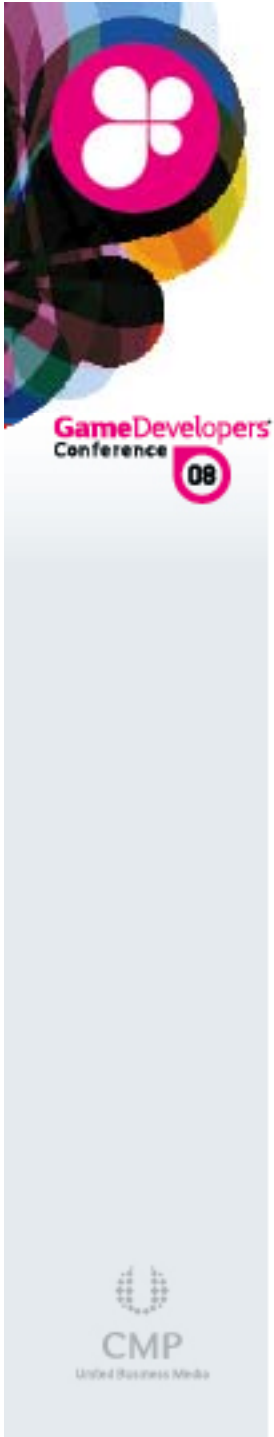
United Business Media

WWW.GDCONF.COM



If ROP is the bottleneck: Solutions

- ④ Render particle effects to lower resolution offscreen texture
 - ④ See GPU Gems 3 chapter by Iain Cantlay
- ④ Disable blending when not needed, especially in larger formats (R32G32B32A32_FLOAT)
- ④ Unbind render targets that are not needed
 - ④ Multiple Render Targets
 - ④ Depth-only passes
- ④ Use R11G11B10 float format for HDR (if you don't need alpha)



If performance is *hitchy* or irregular

- ⌚ Make sure you are not creating/destroying critical resources and shaders at runtime
 - ⌚ Remember to warm caches prior to rendering
- ⌚ Excessive paging when the amount of required video memory is more than available
- ⌚ Could be other engine component like audio, networking, CPU thread synchronization etc.

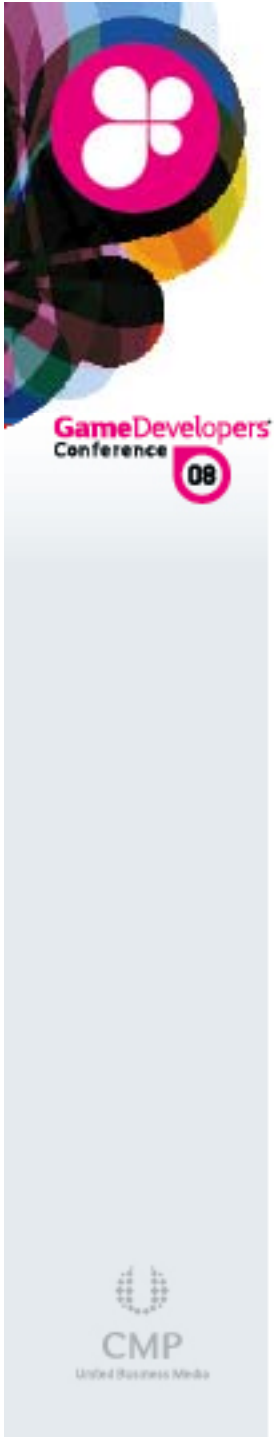


OTHER IHV-SPECIFIC RECOMMENDATIONS



AMD: Clears

- ④ Always clear Z buffer to enable **HiZ**
- ④ Clearing of color render targets is ***not*** free on Radeon HD 2000 and 3000 series
 - Cost is proportional to number of pixels to clear
 - The less pixels to clear the better!
- ④ Here the rule about minimum work applies:
 - Only clear render targets that need to be cleared!
 - Exception for MSAA RTs: need clearing every frame
- ④ RT clears are ***not*** required for optimal multi-GPU usage



AMD: Depth Buffer Formats

- ⊕ Avoid **DXGI_FORMAT_D24_UNORM_S8_UINT** for depth shadow maps
 - Reading back a 24-bit format is a slow path
 - Usually no need for stencil in shadow maps anyway
- ⊕ Recommended depth shadow map formats:
 - DXGI_FORMAT_D16_UNORM**
 - Fastest* shadow map format
 - Precision is enough in most situations
 - ⊕ Just need to set your projection matrix optimally
 - DXGI_FORMAT_D32_FLOAT**
 - High-precision but slower than the 16-bit format



NVIDIA: Clears

- ③ Always Clear Z buffer to enable **ZCULL**
- ③ Always prefer Clears vs. fullscreen quad draw calls
- ③ Avoid partial Clears
 - ③ Note there are no scissored Clears in DX10, they are only possible via draw calls
- ③ Use Clear at the beginning of a frame on any rendertarget or depthstencil buffer
 - ③ In SLI mode driver uses Clears as hint that no inter-frame dependency exist. It can then avoid synchronization and transfer between GPUs



GameDevelopers
Conference 08

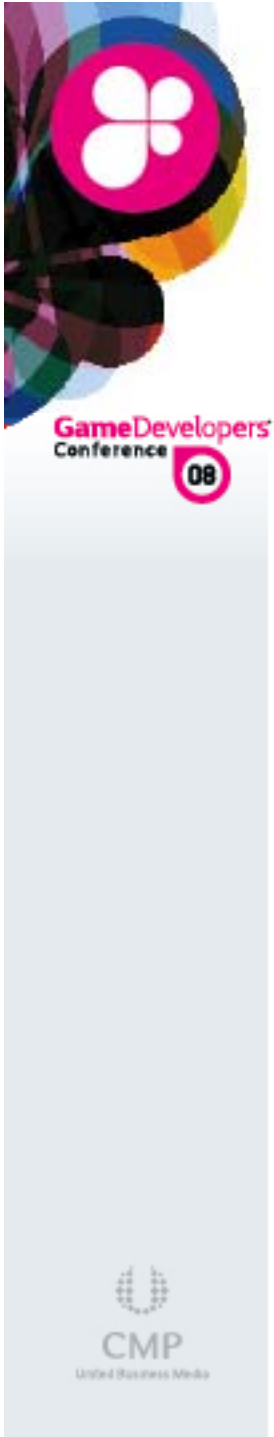
NVIDIA: Depth Buffer Formats

- ④ Use **DXGI_FORMAT_D24_UNORM_S8_UINT**
- ④ **DXGI_FORMAT_D32_FLOAT** should offer very similar performance, but may have lower ZCULL efficiency
- ④ Avoid **DXGI_FORMAT_D16_UNORM**
 - ④ will not save memory or increase performance
- ④ CSAA will increase memory footprint



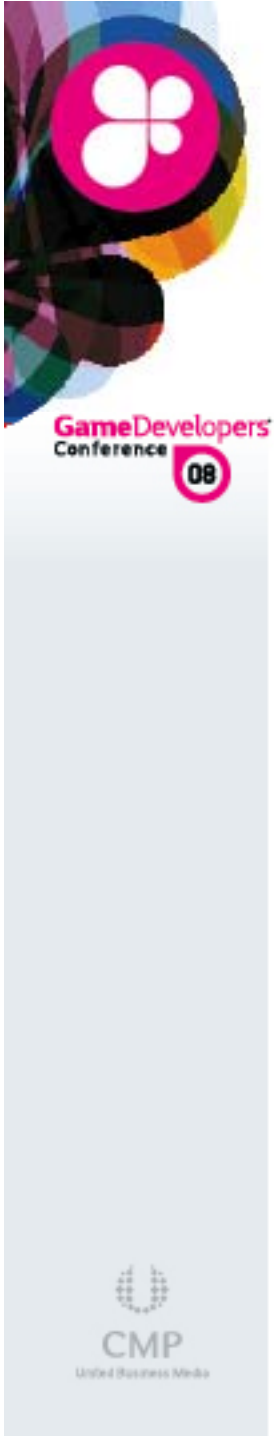
NVIDIA: Attribute Boundedness

- ④ Interleave data when possible into a less VB streams:
 - ④ at least 8 scalars per stream
- ④ Use Load() from Buffer or Texture instead
- ④ Dynamic VBs/IBs might be on system memory accessed over PCIe:
 - ④ maybe CopyResource to USAGE_DEFAULT before using (especially if used multiple times in several passes)
- ④ Passing too many attributes from VS to PS may also be a bottleneck
 - ④ packing and Load() also apply in this case



NVIDIA: ZCULL Considerations

- ⊗ Coarse Z culling is transparent, but it **may underperform** if:
 - ⊗ If depth test changes direction while writing depth (== no Z culling!)
 - ⊗ Depth buffer was written using different depth test direction than the one used for testing (testing is less efficient)
 - ⊗ If stencil writes are enabled while testing (it avoids stencil clear, but may kill performance)
 - ⊗ If DepthStencilView has Texture2D[MS]Array dimension (on GeForce 8 series)
 - ⊗ Using MSAA (less efficient)
 - ⊗ Allocating too many large depth buffers (it's harder for the driver to manage)

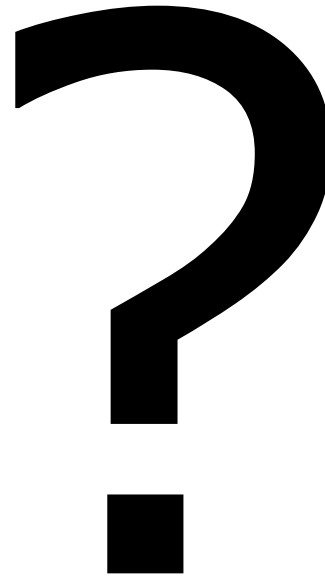


Conclusion

- ④ DX10 is a well-designed and powerful API
- ④ With great power comes great responsibility!
Develop applications with a "DX10" state of mind
A naive port from DX9 will *not* yield expected gains
- ④ Use performance tools available
 - AMD GPUPerfStudio
 - AMD GPUShaderAnalyzer
 - NVIDIA PerfHUD
 - NVIDIA ShaderPerf
- ④ Talk to us



Questions



Ignacio Llamas, NVIDIA
illamas@nvidia.com

Nicolas Thibieroz, AMD
nicolas.thibieroz@amd.com

WWW.GDCONF.COM