

6800 LEAGUES UNDER THE SEA



nVIDIA.®

## Deferred Shading

Shawn Hargreaves

**CLIMAX**™

Mark Harris

NVIDIA





# The Challenge: Real-Time Lighting

- ➊ Modern games use many lights on many objects covering many pixels
  - ➌ computationally expensive
- ➋ Three major options for real-time lighting
  - ➌ Single-pass, multi-light
  - ➌ Multi-pass, multi-light
  - ➌ Deferred Shading
- ➌ Each has associated trade-offs



nVIDIA.



# Comparison: Single-Pass Lighting

For Each Object:

    Render object, apply all lighting in one shader

- ➊ Hidden surfaces can cause wasted shading
- ➋ Hard to manage multi-light situations
  - ➌ Code generation can result in thousands of combinations for a single template shader
- ➌ Hard to integrate with shadows
  - ➍ Stencil = No Go
  - ➎ Shadow Maps = Easy to overflow VRAM



nVIDIA.



# Comparison: Multipass Lighting

For Each Light:

    For Each Object Affected By Light:

```
        framebuffer += brdf( object, light )
```

- ➊ Hidden surfaces can cause wasted shading
- ➋ High Batch Count (1/object/light)
  - ➌ Even higher if shadow-casting
- ➌ Lots of repeated work each pass:
  - ➍ Vertex transform & setup
  - ➎ Anisotropic filtering



nVIDIA.



# Comparison: Deferred Shading

For Each Object:

    Render lighting properties to “G-buffer”

For Each Light:

```
    framebuffer += brdf( G-buffer, light )
```

- ➊ Greatly simplifies batching & engine management
- ➋ Easily integrates with popular shadow techniques
- ➌ “Perfect” O(1) depth complexity for lighting
- ➍ Lots of small lights ~ one big light



nVIDIA.



# Deferred Shading: Not A New Idea!

- ➊ Deferred shading introduced by Michael Deering et al. at SIGGRAPH 1988
  - ➊ Their paper does not ever use the word “deferred”
  - ➊ PixelFlow used it (UNC / HP project)
- ➋ Just now becoming practical for games!

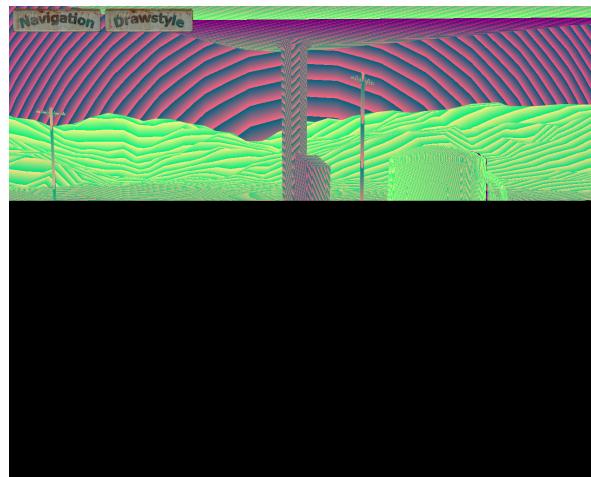
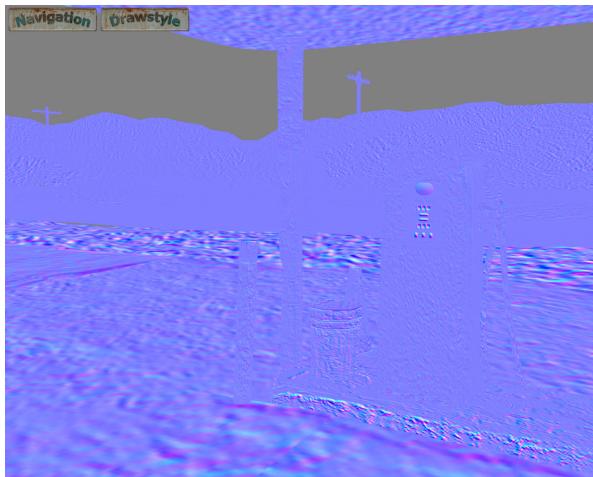


nVIDIA.



# What is a G-Buffer?

- ➊ **G-Buffer = All necessary per-pixel lighting terms**
  - ➊ Normal
  - ➊ Position
  - ➊ Diffuse / Specular Albedo, other attributes
  - ➊ Limits lighting to a small number of parameters!



nVIDIA.



# What You Need

- ➊ Deferred shading is best with high-end GPU features:
  - ➊ Floating-point textures: must store position
  - ➋ Multiple Render Targets (MRT): write all G-buffer attributes in a single pass
  - ➌ Floating-point blending: fast compositing



nVIDIA.



# Attributes Pass

- ➊ Attributes written will depend on your shading
- ➋ Attributes needed
  - ➌ Position
  - ➌ Normal
  - ➌ Color
  - ➌ Others: specular/exponent map, emissive, light map, material ID, etc.
- ➌ Option: trade storage for computation
  - ➍ Store pos.z and compute xy from z + window.xy
  - ➍ Store normal.xy and compute  $z=\sqrt{1-x^2-y^2}$



NVIDIA.



# MRT rules

---

- ➊ Up to 4 active render targets
- ➋ All must have the same number of bits
- ➌ You can mix RTs with different number of channels
- ➍ For example, this is OK:
  - ➎ RT0 = R32f
  - ➎ RT1 = G16R16f
  - ➎ RT2 = ARGB8
- ➎ This won't work:
  - ➏ RT0 = G16R16f
  - ➏ RT1 = A16R16G16B16f



nVIDIA.



# Example MRT Layout

- Three 16-bit Float MRTs

RT1	Diffuse.r	Diffuse.g	Diffuse.b	Specular
RT0	Position.x	Position.y	Position.z	Emissive
RT2	Normal.x	Normal.y	Normal.z	Free

- 16-bit float is overkill for Diffuse reflectance...
  - But we don't have a choice due to MRT rules



nVIDIA.



# Computing Lighting

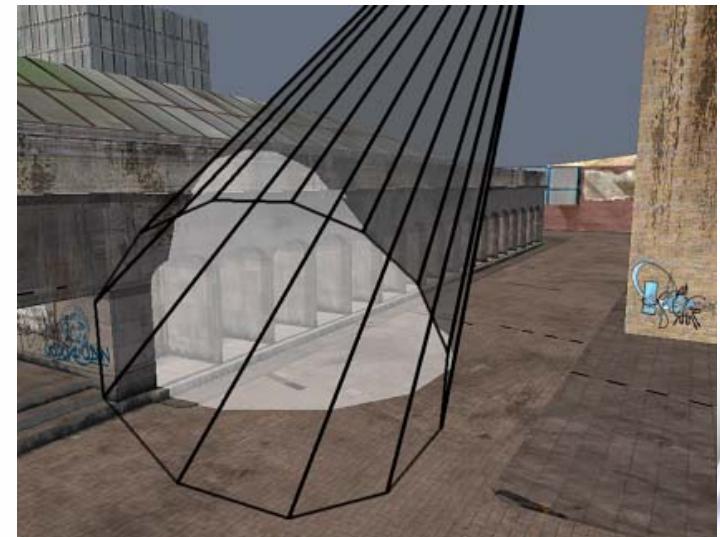
**Render convex bounding geometry**

- Spot Light = Cone
- Point Light = Sphere
- Directional Light = Quad or box

**Read G-Buffer**

**Compute radiance**

**Blend into frame buffer**



Courtesy of Shawn Hargreaves,  
GDC 2004

- Lots of optimizations possible
  - Clipping, occlusion query, Z-cull, stencil cull, etc.



**nVIDIA.**



# Lighting Details

- Blend contribution from each light into accumulation buffer
  - Keep diffuse and specular separate

For each light:

```
diffuse += diffuse(G-buff.N, L)  
specular += G-buff.spec *  
            specular(G-buff.N, G-buff.P, L)
```

- A final full-screen pass modulates diffuse color:

```
framebuffer = diffuse * G-buff.diffuse + specular
```



NVIDIA.



# Options for accumulation buffer(s)

## ➊ Precision

- ➊ 16-bit floating point enables HDR
- ➊ Can use 8-bit for higher performance
  - ➊ Beware of saturation

## ➋ Channels

- ➊ RGBA if monochrome specular is enough
- ➊ 2 RGBA buffers if RGB diffuse and specular are both needed.
- ➊ Small shader overhead for each RT written

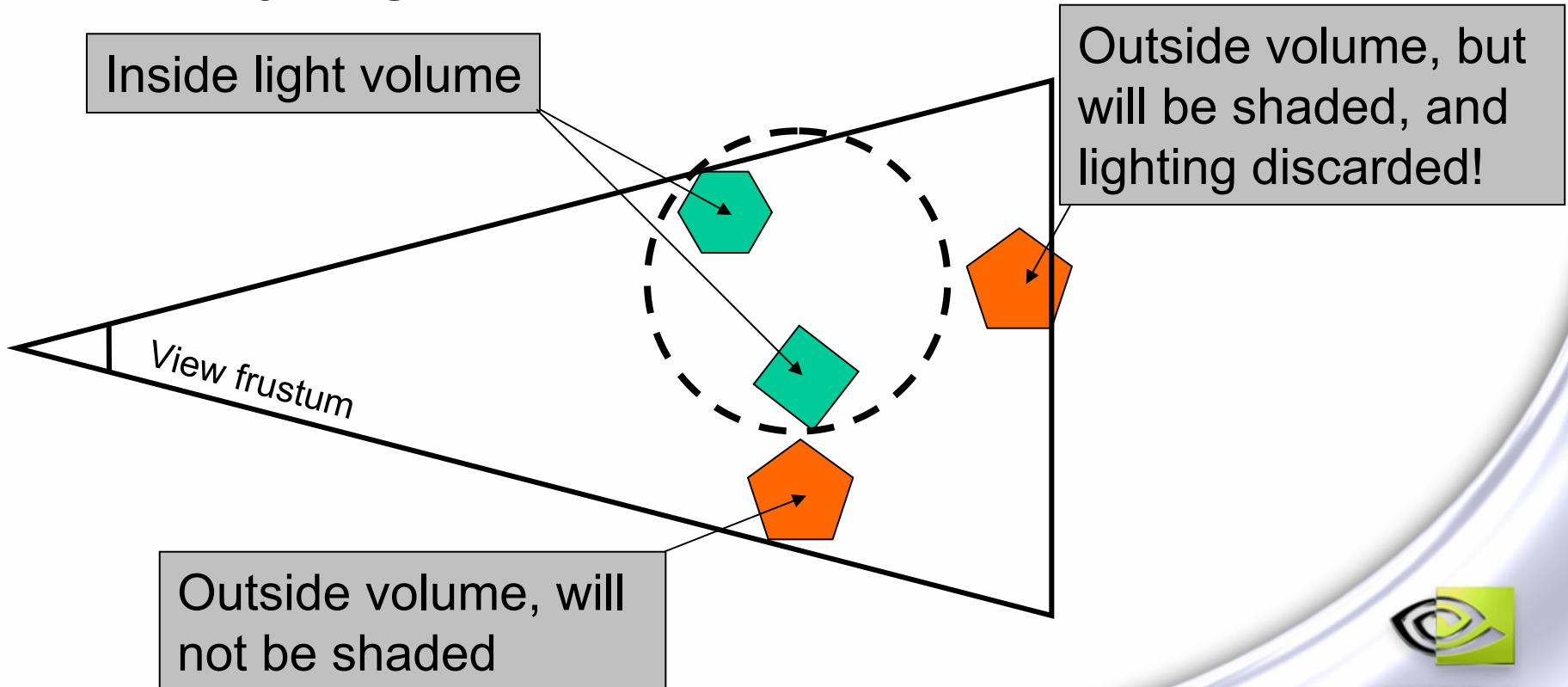


nVIDIA.



# Lighting Optimization

- Only want to shade surfaces inside light volume
  - Anything else is wasted work



nVIDIA.



# Optimization: Stencil Cull

- ➊ Two pass algorithm, but first pass is very cheap
  - ➌ Rendering without color writes = 2x pixels per clock

## 1. Render light volume with color write disabled

- ➌ Depth Func = LESS, Stencil Func = ALWAYS
- ➌ Stencil Z-FAIL = REPLACE (with value X)
- ➌ Rest of stencil ops set to KEEP

## 2. Render with lighting shader

- ➌ Depth Func = ALWAY, Stencil Func = EQUAL,  
all ops = KEEP, Stencil Ref = X
- ➌ Unlit pixels will be culled because stencil will not  
match the reference value

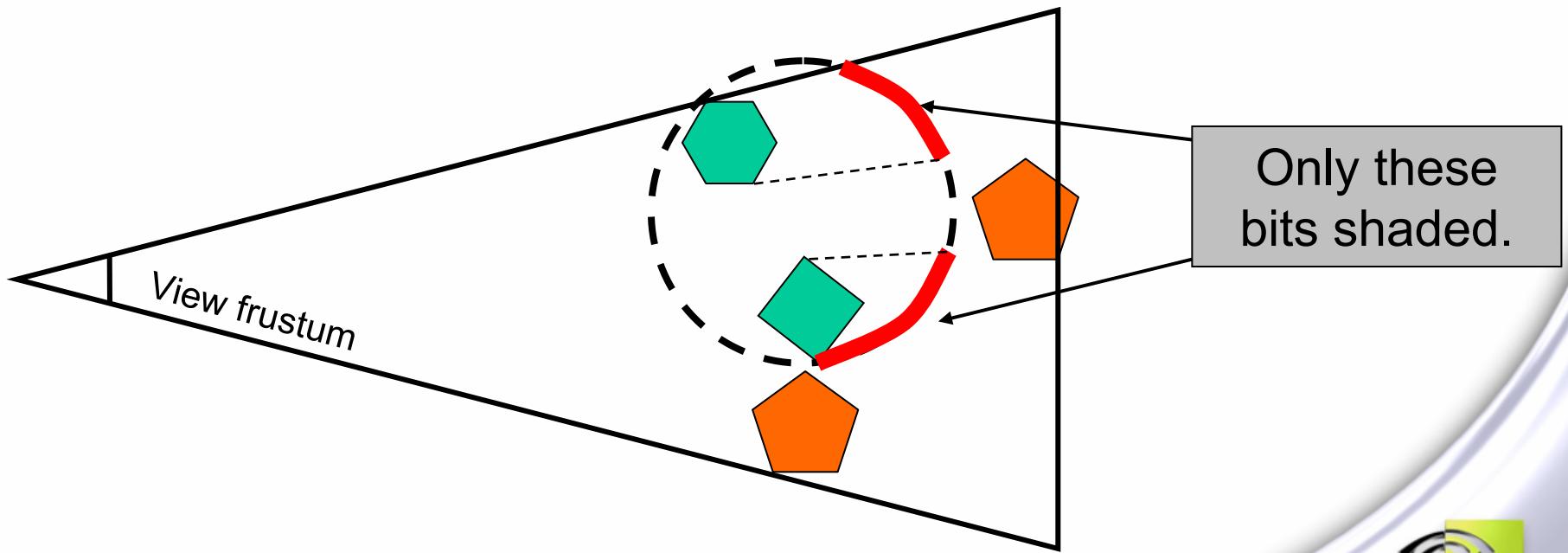


nVIDIA.



# Setting up Stencil Buffer

- Only regions that fail depth test represent objects within the light volume



nVIDIA.



# Shadows

- ➊ **Shadow maps work very well with deferred shading**
  - ➌ Work trivially for directional and spot lights
  - ➌ Point (omni) lights are trickier...
- ➋ **Don't forget to use NVIDIA hardware shadow maps**
  - ➌ Render to shadow map at 2x pixels per clock
  - ➌ Shadow depth comparison in hardware
  - ➌ 4 sample percentage closer filtering in hardware
  - ➌ Very fast high-quality shadows!
- ➌ May want to increase shadow bias based on pos.z
  - ➌ If using fp16 for G-buffer positions



nVIDIA.



# Virtual Shadow Depth Cube Texture

- ➊ Solution for point light shadows
  - ➊ Technique created by Will Newhall & Gary King
- ➋ Unrolls a shadow cube map into a 2D depth texture
  - ➊ Pixel shader computes ST and depth from XYZ
  - ➋ G16R16 cubemap efficiently maps XYZ->ST
  - ➌ Free bilinear filtering offsets extra per-pixel work
- ➌ More details in *ShaderX<sup>3</sup>*
  - ➊ Charles River Media, October 2004



nVIDIA.



# Multiple Materials w/ Deferred Shading

- ➊ Deferred shading doesn't scale to multiple materials
  - ➌ Limited number of terms in G-buffer
  - ➌ Shader is tied to light source – 1 BRDF to rule them all
  
- ➋ Options:
  - ➌ Re-render light multiple times, 1 for each BRDF
    - ➌ Loses much of deferred shading's benefit
  - ➌ Store multiple BRDFs in light shader, choose per-pixel
    - ➌ Use that last free channel in G-buffer to store material ID
    - ➌ Reasonably coherent dynamic branching
    - ➌ Should work well on pixel shader 3.0 hardware



nVIDIA.



# Transparency

- ➊ Deferred shading does not support transparency
  - ➌ Only shades nearest surfaces
- ➋ Just draw transparent objects last
  - ➌ Can use depth peeling
  - ➌ Blend into final image, sort back-to-front as always
  - ➌ Use “normal” shading / lighting
  - ➌ Make sure you use the same depth buffer as the rest
- ➌ Also draw particles and other blended effects last



nVIDIA.



# Post-Processing

- ➊ G-buffer + accum buffers can be used as input to many post-process effects
  - ➊ Glow
  - ➊ Auto-Exposure
  - ➊ Distortion
  - ➊ Edge-smoothing
  - ➊ Fog
  - ➊ Whatever else!
  - ➊ HDR
- ➋ See HDR talk



nVIDIA.



# Anti-Aliasing with Deferred Shading

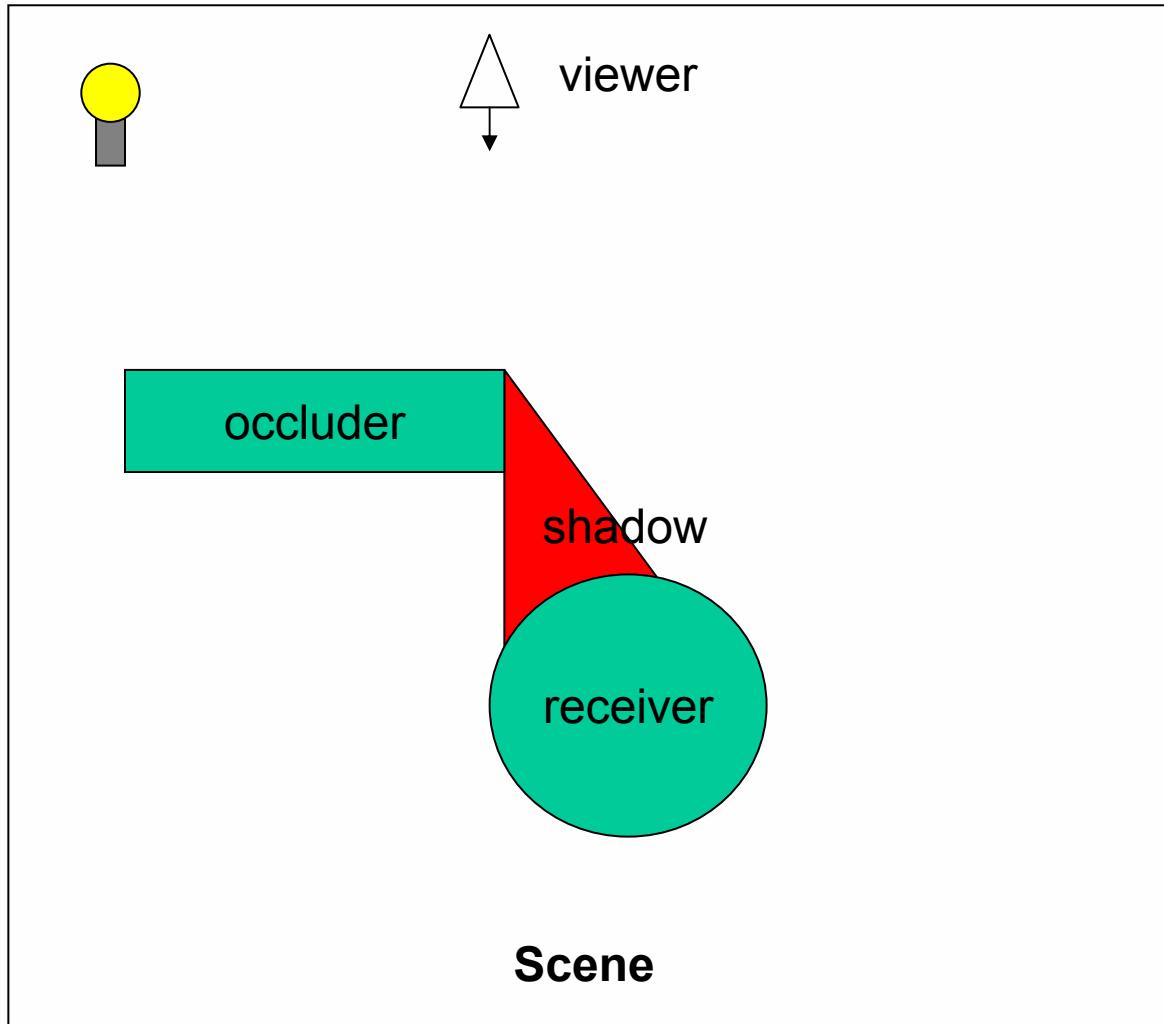
- ➊ Deferred shading is incompatible with MSAA
- ➋ API doesn't allow antialiased MRTs
  - ➌ But this is a small problem...
- ➌ AA resolve has to happen after accumulation!
  - ➍ Resolve = process of combining multiple samples
- ➍ G-Buffer cannot be resolved
  - ➎ What happens to an FP16 position when resolved?



nVIDIA.



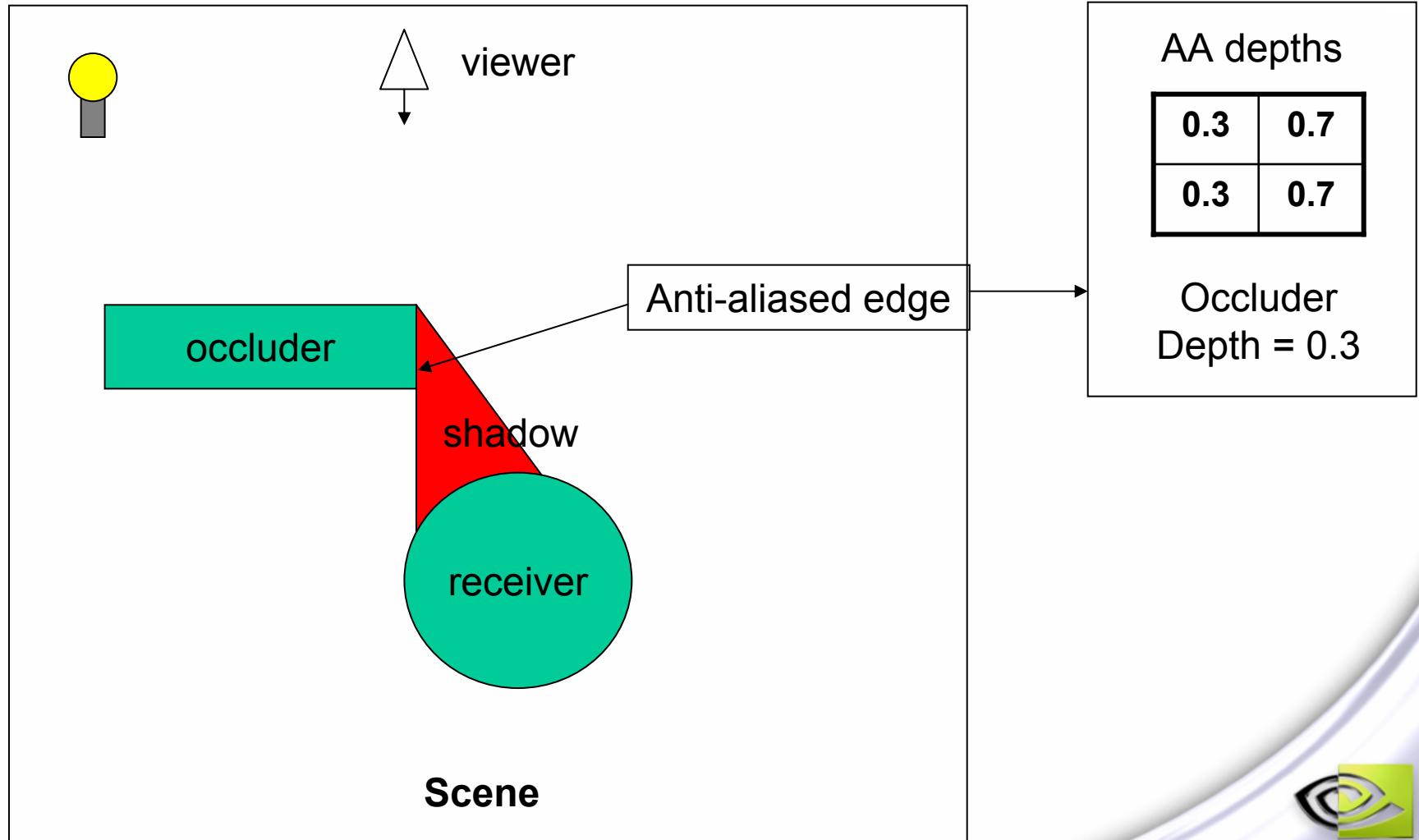
# Shadow Edge, Correct AA Resolve



nVIDIA.



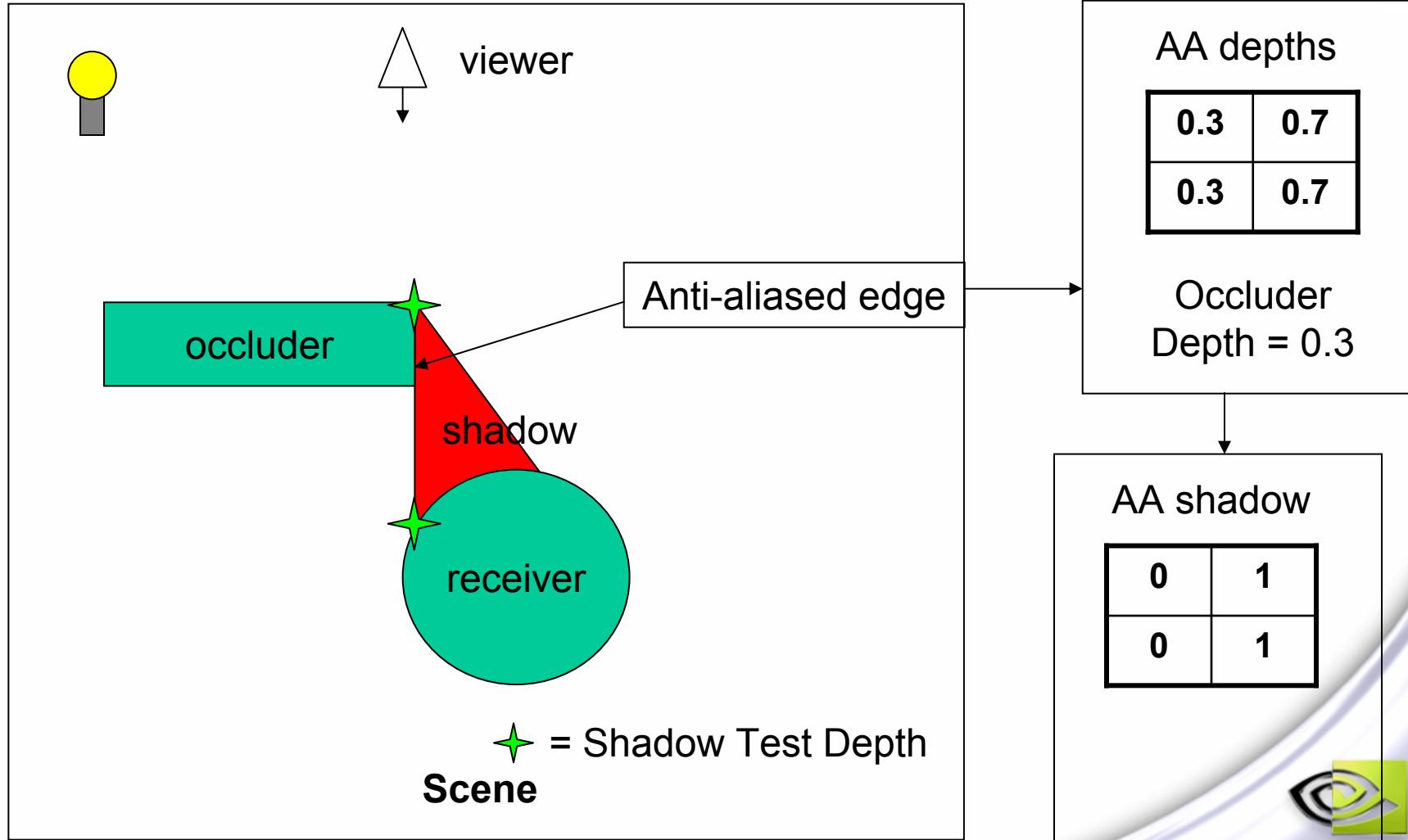
# Shadow Edge, Correct AA Resolve



nVIDIA.



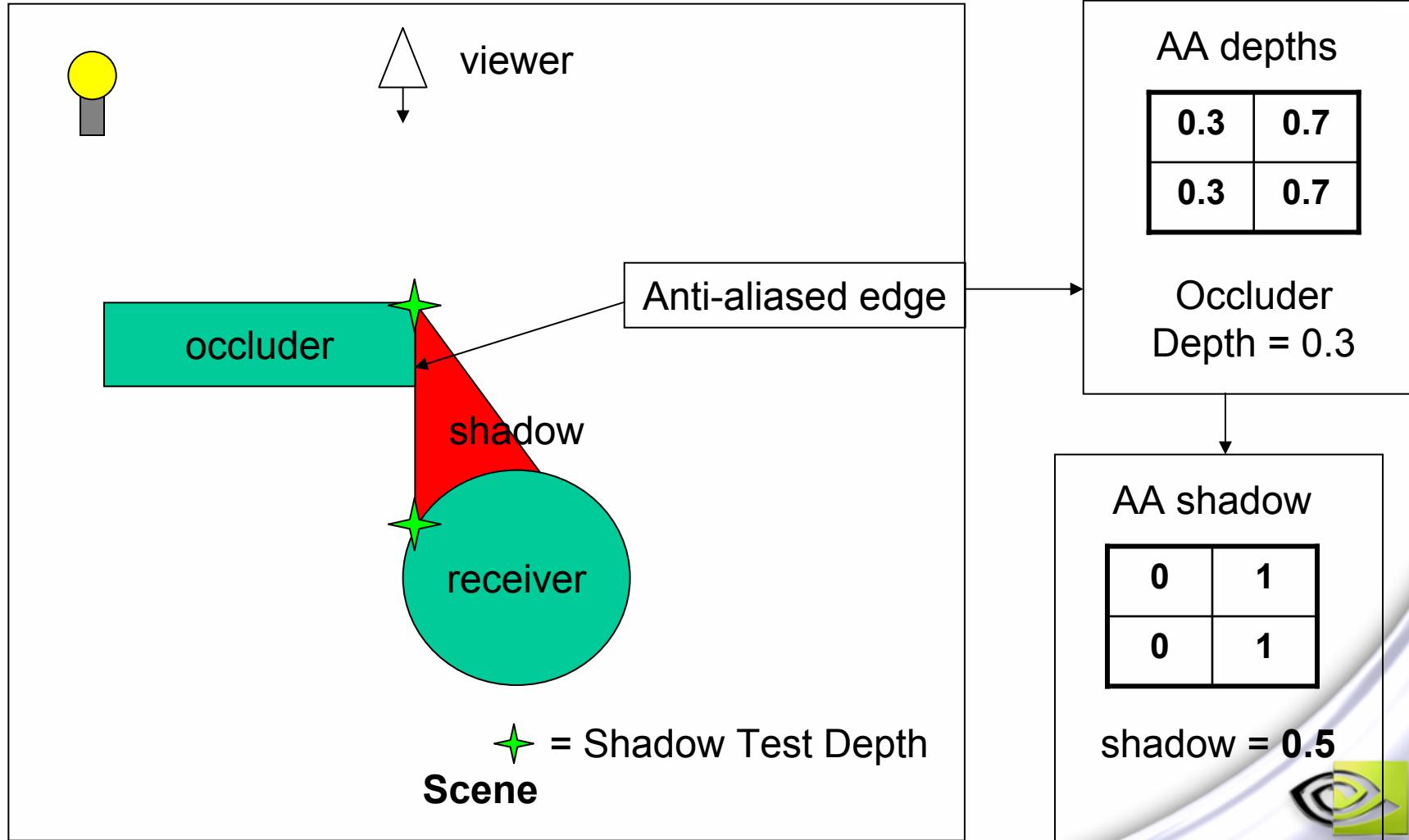
# Shadow Edge, Correct AA Resolve



nVIDIA.



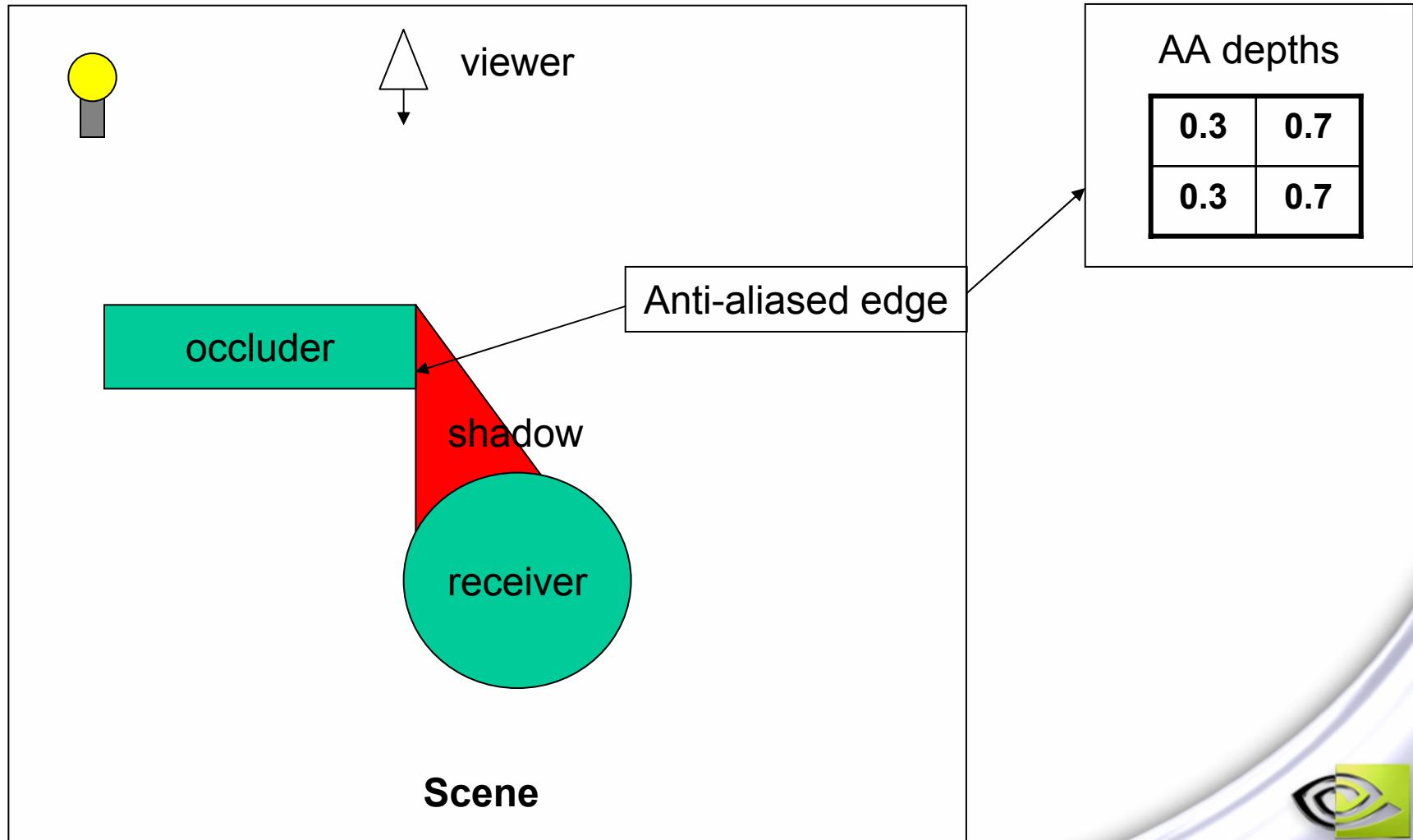
# Shadow Edge, Correct AA Resolve



NVIDIA.



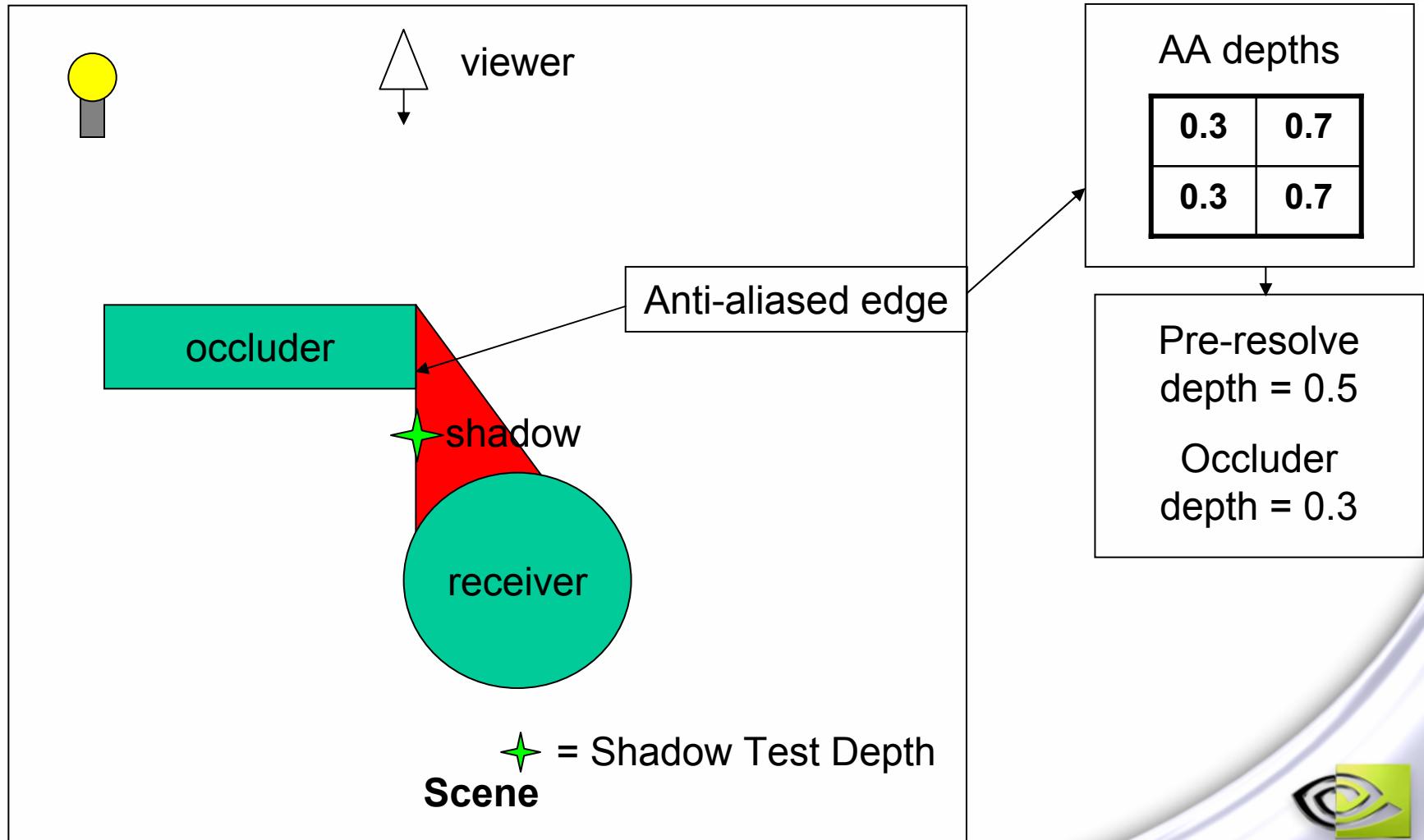
# Shadow Edge, G-Buffer Resolve



nVIDIA.



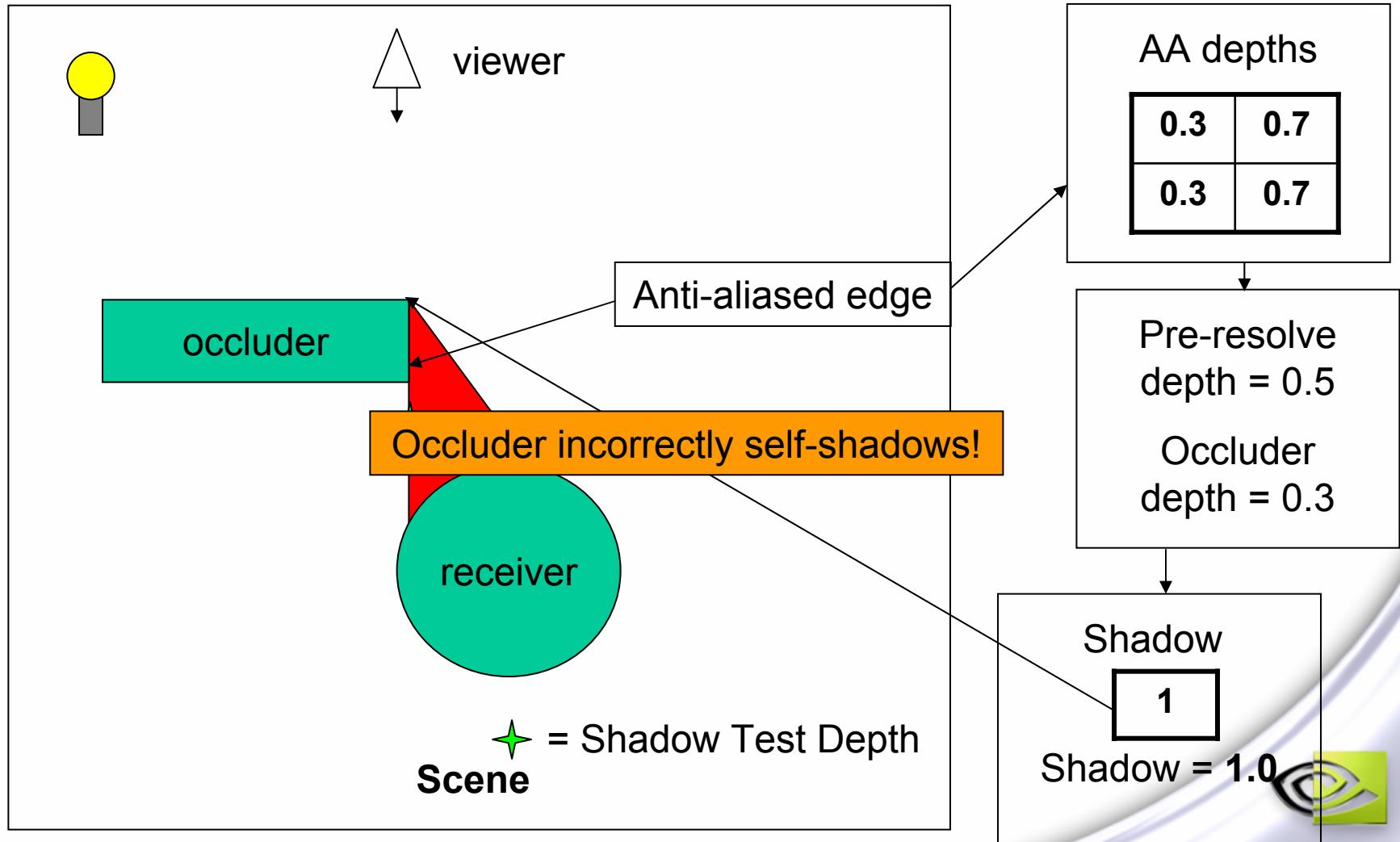
# Shadow Edge, G-Buffer Resolve



nVIDIA.



# Shadow Edge, G-Buffer Resolve





# Other AA options?

- ➊ **Supersampling lighting is a costly option**
  - ➌ Lighting is typically the bottleneck, pixel shader bound
  - ➌ 4x supersampled lighting would be a big perf. Hit
- ➋ **“Intelligent Blur” : Only filter object edges**
  - ➌ Based on depths and normals of neighboring pixels
  - ➌ Set “barrier” high, to avoid interior blurring
  - ➌ Full-screen shader, but cheaper than SSAA



nVIDIA.



# Should I use Deferred Shading?

- ➊ This is an **ESSENTIAL** question
- ➋ Deferred shading is not always a win
  - ➌ One major title has already scrapped it!
  - ➌ Another came close
- ➌ Many tradeoffs
  - ➍ AA is problematic
  - ➍ Some scenes work well, others very poorly
- ➌ The benefit will depend on your application
  - ➍ Game design
  - ➍ Level design



nVIDIA.



# When is Deferred Shading A Win?

- ➊ Not when you have many directional lights
  - ➌ Shading complexity will be  $O(R*L)$ ,  $R$  = screen res.
  - ➌ Outdoor daytime scenes probably not a good case
- ➋ Better when you have lots of local lights
  - ➌ Ideal case is non-overlapping lights
  - ➌ Shading complexity  $O(R)$
  - ➌ Nighttime scenes with many dynamic lights!
- ➌ In any case, make sure G-Buffer pass is cheap



nVIDIA.



# Gosh, what about z-cull & SM3.0?

- ➊ Isn't the goal of z-cull to achieve deferred shading?
  - ➊ Do an initial front-to-back-sorted z-only pass.
  - ➊ Then you will shade only visible surfaces anyway!
- ➋ Shader Model 3.0 allows “uber shaders”
  - ➊ Iterate over multiple lights of different types in “traditional” (non-deferred) shading
- ➌ Combine these, and performance could be as good (better?) than deferred shading!
  - ➊ More tests needed



nVIDIA.



# We don't have all the answers

- ➊ We can't tell you to use it or not
  - ➊ Experimentation and analysis is important
  - ➊ Depends on your application
  - ➊ Need to have a fallback anyway



nVIDIA.



**Sorry to end it this way, but...**

---

**MORE RESEARCH IS NEEDED!  
PLEASE SHARE YOUR FINDINGS!**

**(you can bet we'll share ours)**



**nVIDIA.**



# Questions?

- <http://developer.nvidia.com>
- [mharris@nvidia.com](mailto:mharris@nvidia.com)



nVIDIA®



# GeForce 6800 Guidance (1 of 6)

- ➊ Allocate render targets FIRST
  - ➌ Deferred Shading uses many RTs
  - ➌ Allocating them first ensures they are in fastest RAM
  
- ➋ Keep MRT usage to 3 or fewer render targets
  - ➌ Performance cliff at 4 on GeForce 6800
  - ➌ Each additional RT adds shader overhead
  - ➌ Don't render to all RTs if surface doesn't need them
    - ➌ e.g. Sky Dome doesn't need normals or position



nVIDIA.



# GeForce 6800 Guidance (2 of 6)

- ➊ Use aniso filtering during G-buffer pass
  - ➌ Will help image quality on parts of image that don't benefit from "edge smoothing AA"
  - ➌ Only on textures that need it!
  
- ➋ Take advantage of early Z- and Stencil culling
  - ➌ Don't switch z-test direction mid-frame
  - ➌ Avoid frequent stencil reference / op changes



nVIDIA.



# GeForce 6800 Guidance (3 of 6)

- ➊ Use hardware shadow mapping (“UltraShadow”)
  - ➌ Use D16 or D24X8 format for shadow maps
  - ➌ Bind 8-bit color RT, disable color writes on updates
  - ➌ Use tex2Dproj to get hardware shadow comparison
  - ➌ Enable bilinear filtering to get 4-sample PCF



nVIDIA.



# GeForce 6800 Guidance (4 of 6)

- ➊ Use fp16 filtering and blending
  - ➌ Fp16 textures are fully orthogonal!
  - ➌ No need to “ping-pong” to accumulate light sources
- ➋ Use the lowest precision possible
  - ➌ Lower-precision textures improve cache coherence, reduce bandwidth
  - ➌ Use half data type in shaders



nVIDIA.



# GeForce 6800 Guidance (5 of 6)

- ➊ Use write masks to tell optimizer sizes of operands
  - ➊ Can schedule multiple instructions per cycle
    - ➊ Two simultaneous 2-component ops, or
    - ➋ One 3-component op + 1 scalar op
- ➋ Without write masks, compiler must be conservative



nVIDIA.



# GeForce 6800 Guidance (6 of 6)

- Use fp16 normalize()
  - Compiles to single-cycle nrmh instruction
  - Only applies to half3, so:

```
half3 n = normalize(tex2D(normalmap, coords).xyz); // fast  
half4 n = normalize(tex2D(normalmap, coords)); // slow  
float3 n = normalize(tex2D(normalmap, coords).xyz); // slow
```



NVIDIA.



# Example Attribute Layout

- ➊ **Normal: x,y,z**
- ➋ **Position: x, y, z**
- ➌ **Diffuse Reflectance: RGB**
- ➍ **Specular Reflectance (“Gloss Map”, single channel)**
- ➎ **Emissive (single channel)**
- ➏ **One free channel**
  - ➐ **Ideas on this later**
  - ➑ **Your application will dictate**



nVIDIA.