

# learn network inspire

GameDevelopers'  
Conference

08



February 18-22, 2008  
San Francisco

[www.gdconf.com](http://www.gdconf.com)



GameDevelopers  
Conference

08

# Practical Spherical Harmonics Based PRT Methods

**Manny Ko**

Naughty Dog

**Jerome Ko**

UCSD / Bunkspeed



CMP

United Business Media

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



GameDevelopers  
Conference 08

# Outline

- Background ambient occlusion, HL2
- n Spherical Harmonics theory
- i Pre-processing
  - PRT compression – 4 to 6 bytes per sample
- Conclusion and game scene demo



# Background

GameDevelopers  
Conference 08

- Many variants of Precomputed Radiance Transfer – self-shadow, interreflections, diffuse vs. glossy
- e Power of PRT best understood in the context of ambient occlusion (AO) and HL2 basis



# Goals

- Diffuse *self-shadowing* for rigid bodies (aka neighborhood transfer)
- n Generalizes easily to interreflections
- l Decouples *visibility* calculation from lighting
- m Pre-integrate *surface variations*



# Ambient Occlusion

GameDevelopers  
Conference

08

- n Pioneered by ILM  
Accessibility shading



CMP

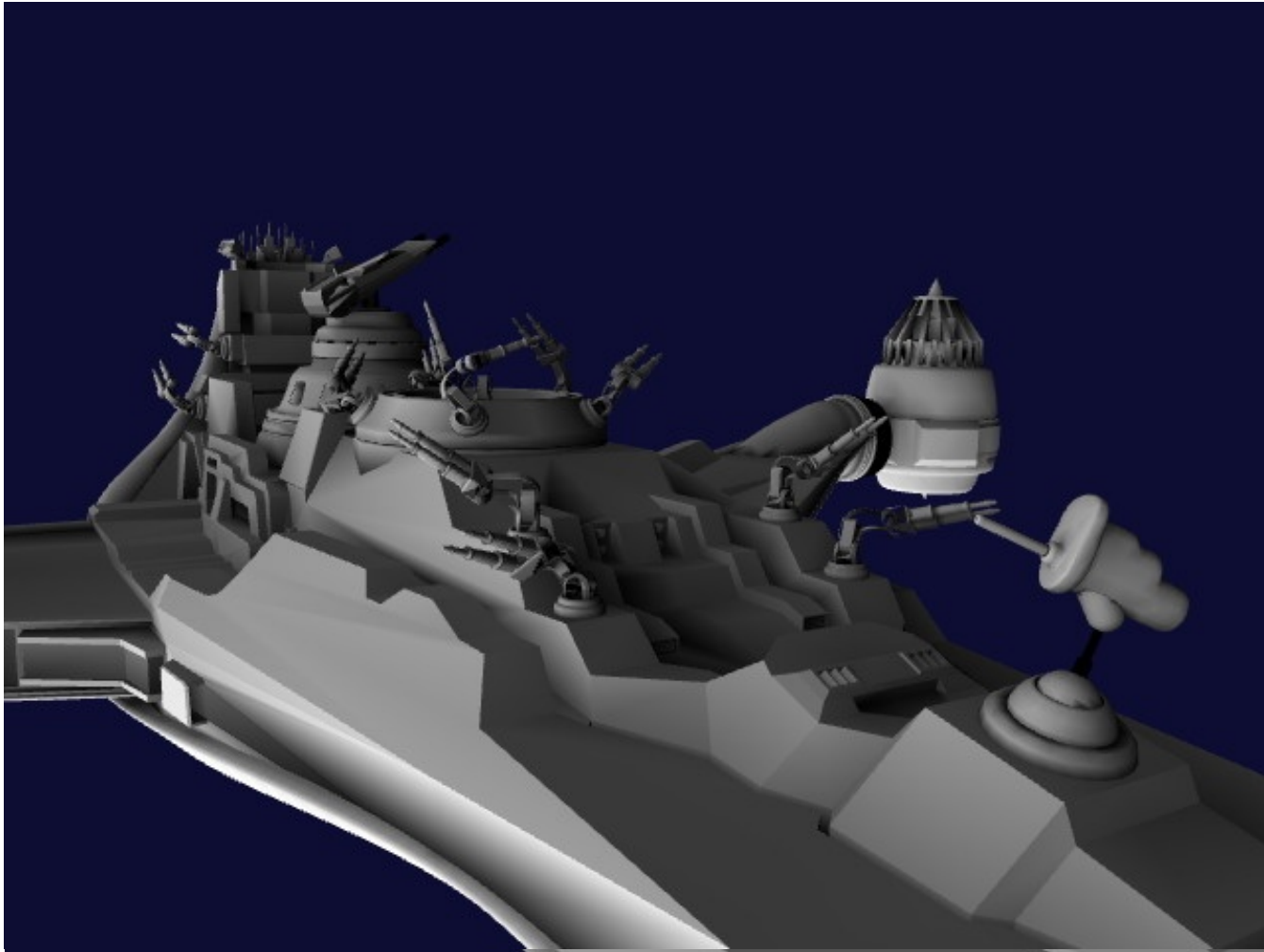
United Business Media

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



# Ambient Occlusion

Game Developers  
Conference  
08





# Ambient Occlusion Flavors

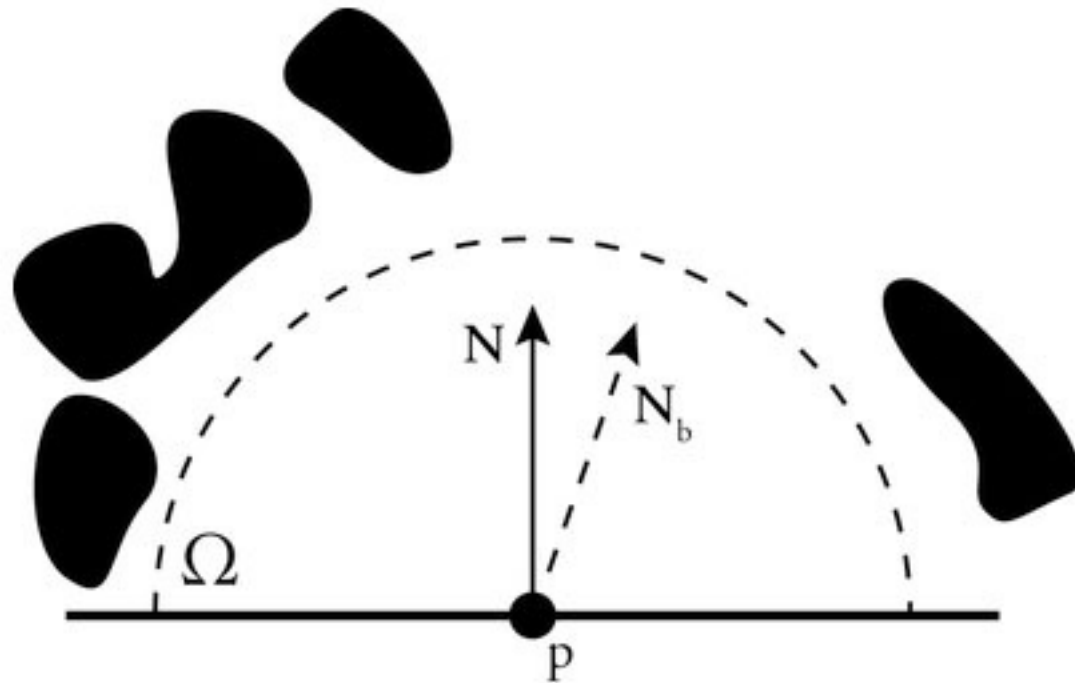
GameDevelopers  
Conference 08

Single scalar – cheap but no directional response

c Bent Normal



# Bent Normal: basic idea



$$A_p = \frac{1}{\pi} \int_{\Omega} V_{p,\omega}(N \cdot \omega) d\omega$$



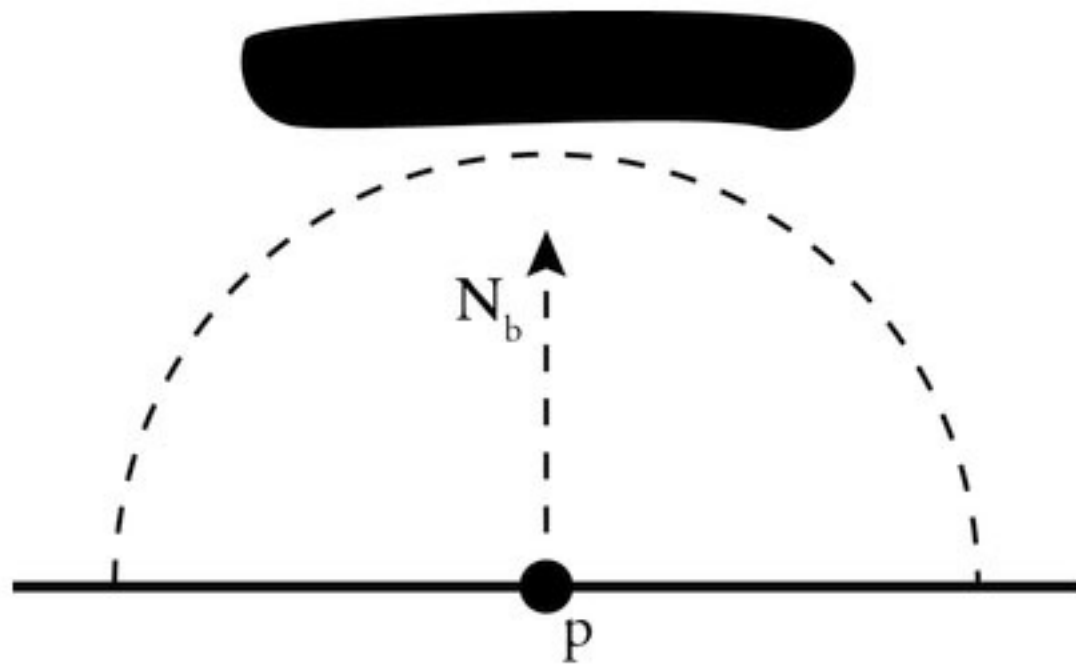
# Bent Normal

Half3 + scalar AO = 7 bytes per sample

- Gives some directional variations but it does not always work



# Bent Normal: fail case





# Bent Normal: next step

- t Overcome the single direction, single weight limitation

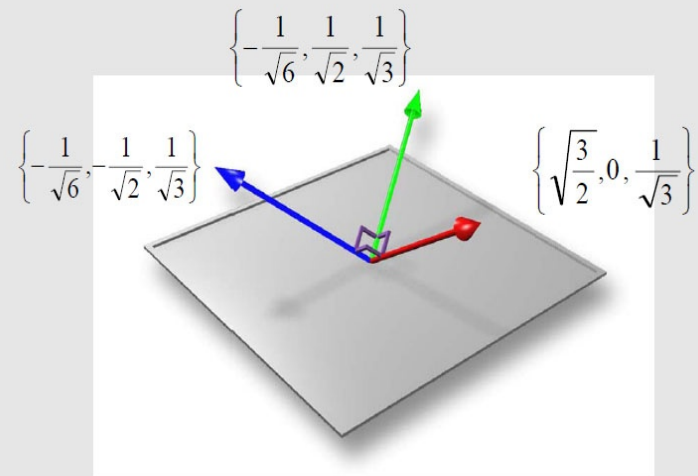


# HL2 Basis

Valve's "Radiosity Normal Mapping"

3 RGB basis vector per sample

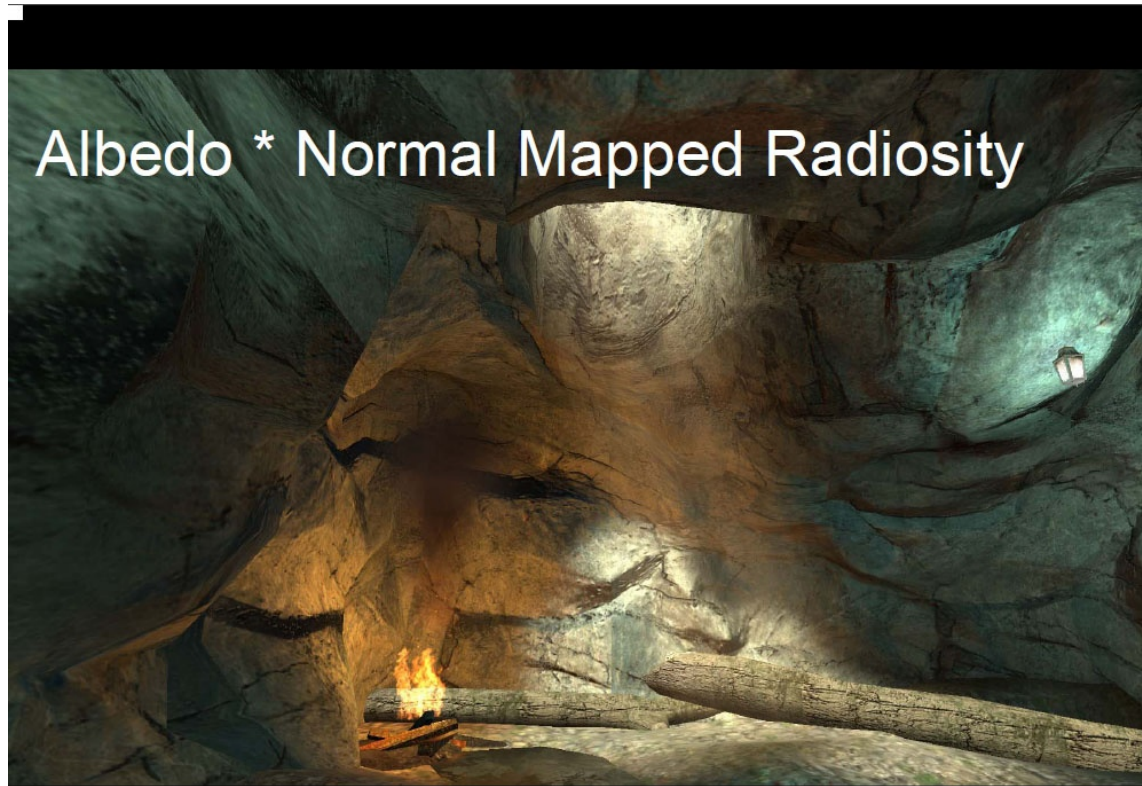
Basis for Radiosity Normal Mapping





# HL2 Samples

Game Developers  
Conference 08





# HL2 Strengths

- e Successfully integrated GI with surface-details
- r Surface details respond to runtime lights
- s 1<sup>st</sup> attempt in games to use a 'lighting basis'
- s Orthonormal basis
- s Good for normal-maps



# HL2 Weakness

- Lots of storage – 3 RGB maps per sample plus UVs
- 3 texture lookup + 3 dotp/ pixel after rotating basis into tangent space. Shader cost can be a problem.
- h Incomplete coverage of the hemisphere,  
3 directions not enough?



# Key Ideas for PRT

Game Developers  
Conference 08

- T Visibility is the key bottleneck since it involves sampling the scene
- n Take advantage of off-line compute
- f Separate pre-integration of visibility from that of incoming light
- e Find a basis to store these functions



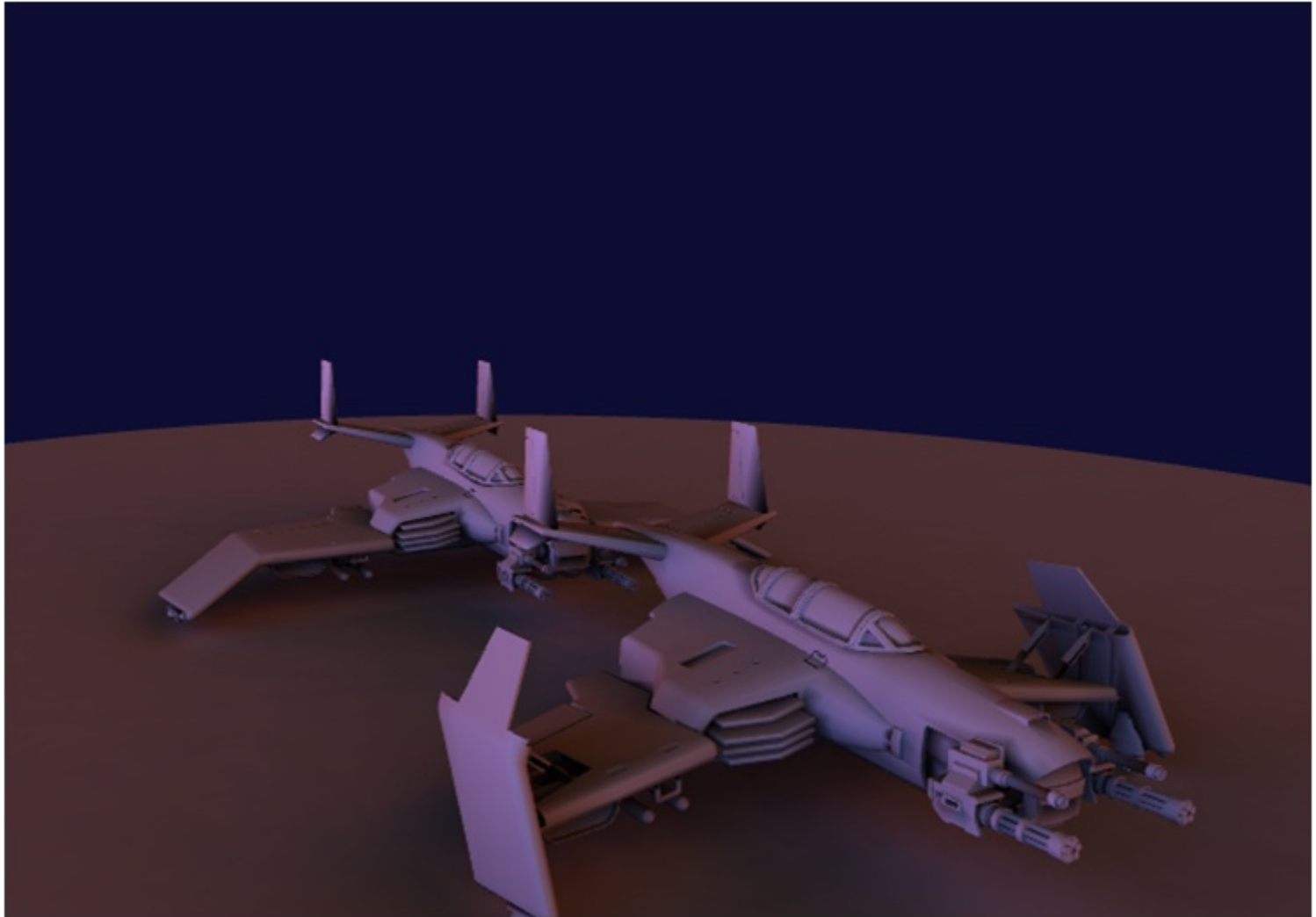
# Key Ideas for PRT

- T SH is a good starting point for a basis
- r We do not trace rays using the GPU! Nor do we take tons of samples inside the pixel shader
- e Better quality than AO or HL2 but more efficient than latter
- h Works well for static scene + moving lights



# Hovercraft from Warhawk

GameDevelopers  
Conference 08





# Demo

Buddha – 1.08 m triangles

33 fps on ATI X700

250 fps on G8800 GTS

Vertex shader only method

Vertex buffer not optimized, could be faster

t Dragon – 871k triangles

40.7 fps on ATI X700

# Implementation

Jerome Ko

UCSD / Bunkspeed



# Outline

- ⌘ Short introduction to SH (very short!)
- o Environment map projection
- p PRT coefficient generation
- g Runtime reconstruction in vertex shader



# What are we trying to solve?

$$E = \int_S Li(\mathbf{x}) \cdot BRDF(\mathbf{x}, N) \cdot V(p, \mathbf{x}) dx$$

- n Use SH to simplify each of these functions



# Spherical Harmonics

GameDevelopers  
Conference 08

- i 2D Fourier series on the sphere
- s Represent 2D signals as scaled and shifted Fourier basis
- n Perfect for spherical functions



# Basis Function Projection

GameDeveloper:  
Conference 08

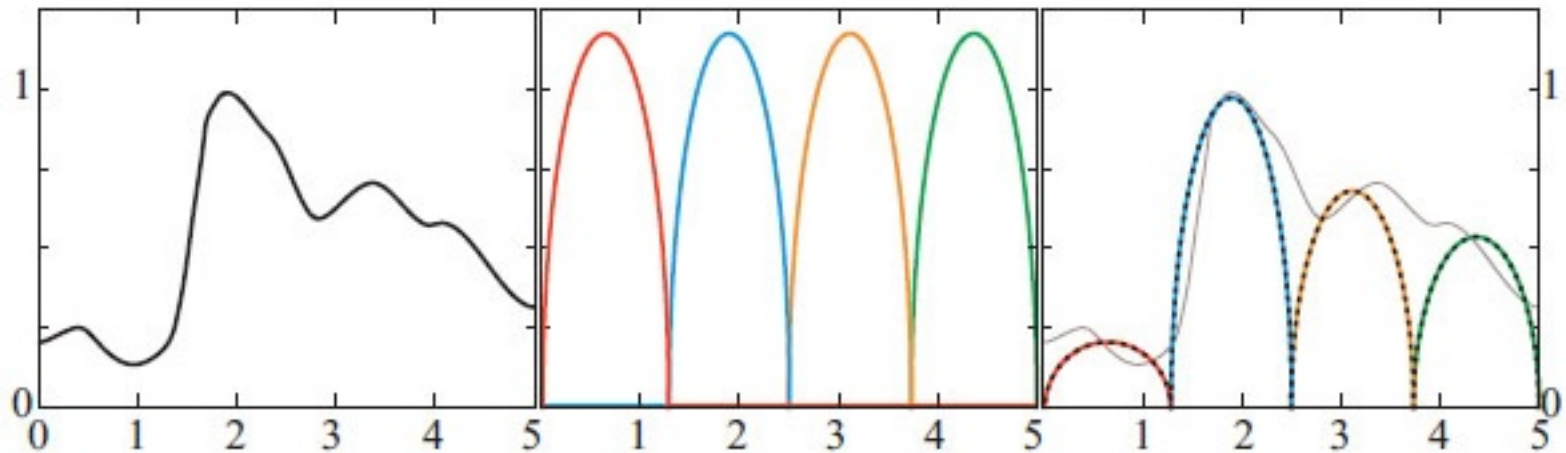


Image courtesy of Real Time Rendering 3<sup>rd</sup> Edition



# Spherical Harmonics Properties

- i Orthonormal basis ->
  - convolution as dot products
  - simple projection
- n Multi-resolution/band-limited, related to Fourier decomposition
- / Solid math foundation
- a Stable rotation, no wobbling lights
  - Can add/subtract, *lerp*.



# SH: bands ( $l, m$ )

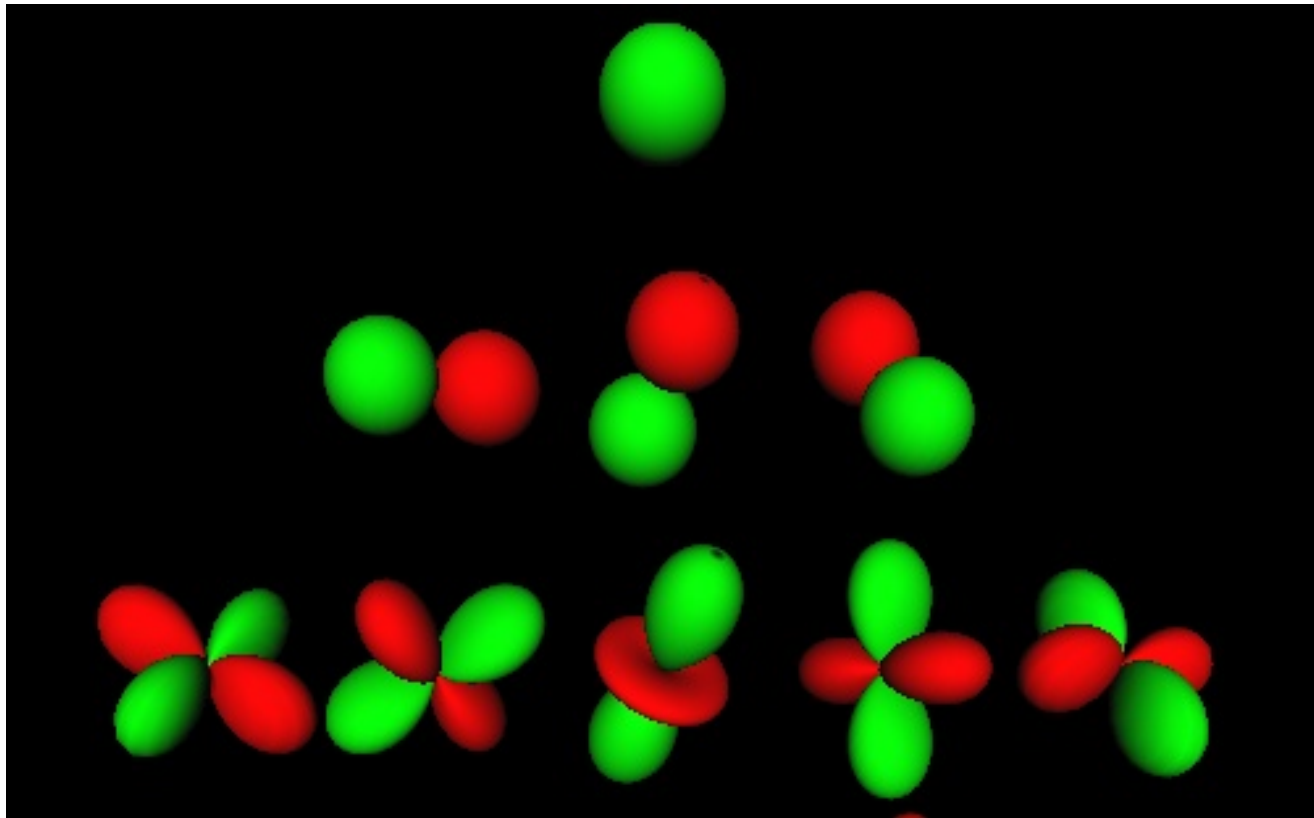
SH-order: order is the number of bands. ' $l$ ' is the order  
0<sup>th</sup> order: just a constant (DC) which is your AO term

- 1<sup>st</sup> order: 3 linear terms which is your 'bent normal' term
- 2<sup>nd</sup> order: 5 terms – these give the extra directional response we want
- In general ( $2l+1$ ) terms



# Spherical Harmonic Functions

Game Developers  
Conference  
08





# Spherical Harmonics - Issues

Game Developers  
Conference 08

- i Scary math!
- i Looks difficult to implement
- t Lots of coefficients, too much storage?
- e Difficult or time consuming to generate?



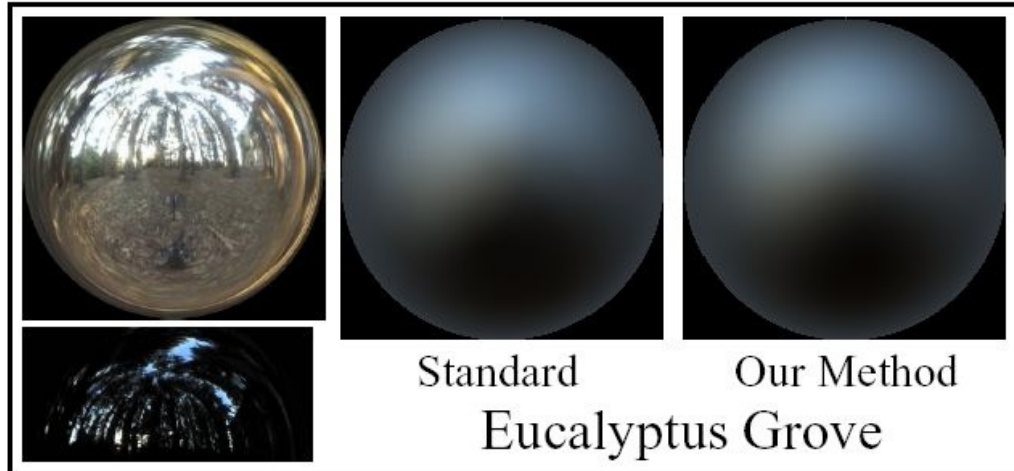
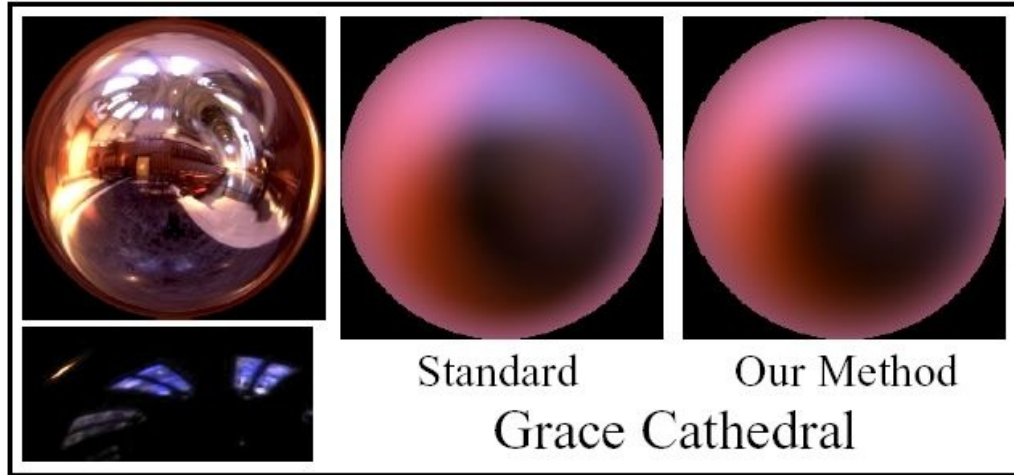
# Environment Map Projection

- P First we need a light source
- I Ravi R. introduced a technique to light diffuse surfaces with distant environment map illumination
- S Code available on Ravi's site, we integrated it into our demo app



# Environment Map Projection

GameDevelopers  
Conference 08





# Environment Map Projection

- P 9 SH-compressed coefficients
- c Instead of a large cubemap you only need 27 numbers!
- g Secret lies in the solid angle formula
- h For each pixel evaluate the SH basis, then weight it by the solid angle



# Solid Angle

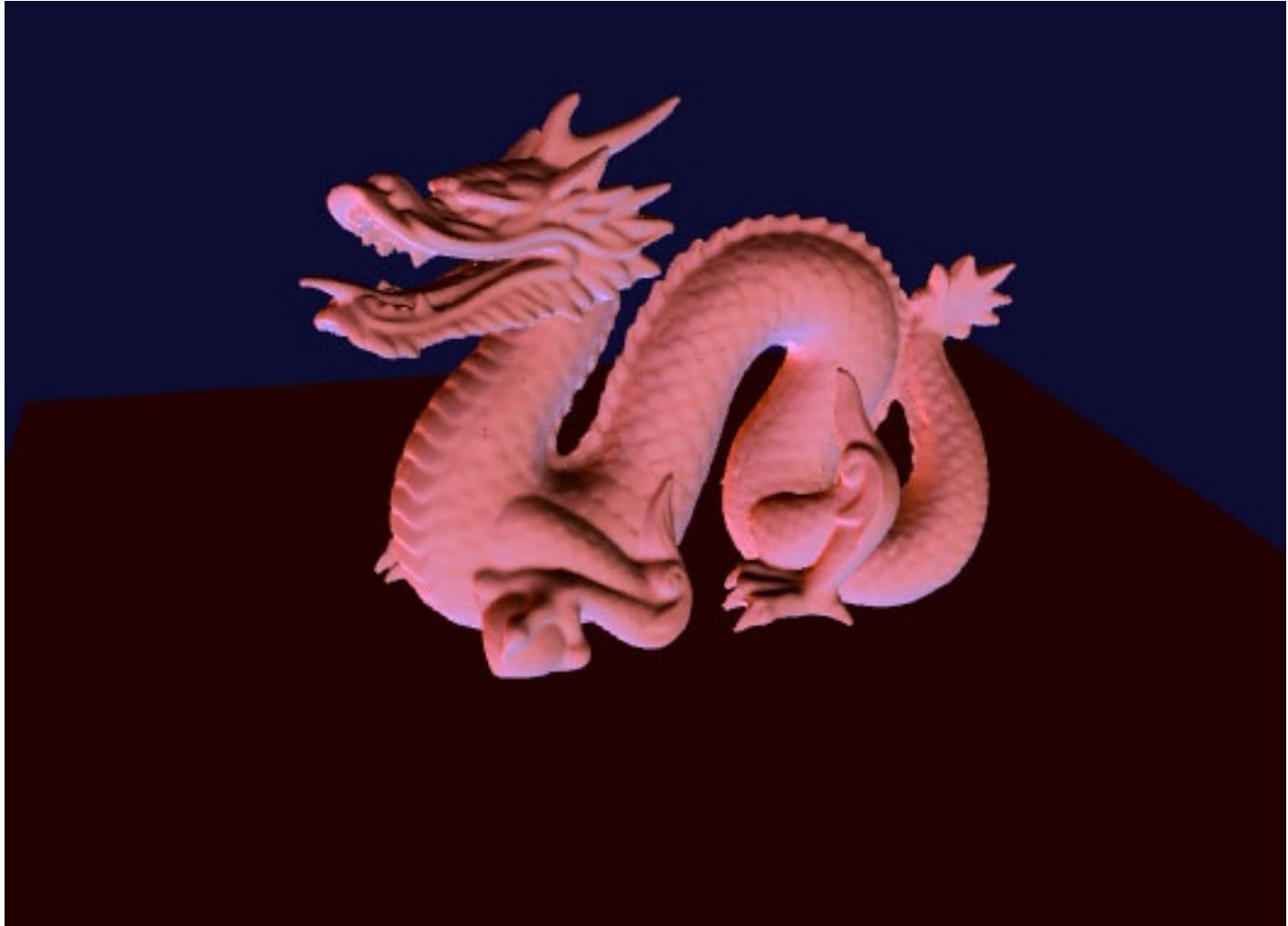
GameDevelopers  
Conference 08

```
double omega = 2*PI/width *  
              2*PI/height *  
              sinc(theta);
```



# Envmap Demo

Game Developers  
Conference 08





# PRT Coefficient Generation

Game Developers  
Conference 08

**G** We have simplified the light source, now we must also project the visibility function using SH

This gives us compact storage of self-shadowing information

**m** Encodes the directional variations more accurately than AO or HL2



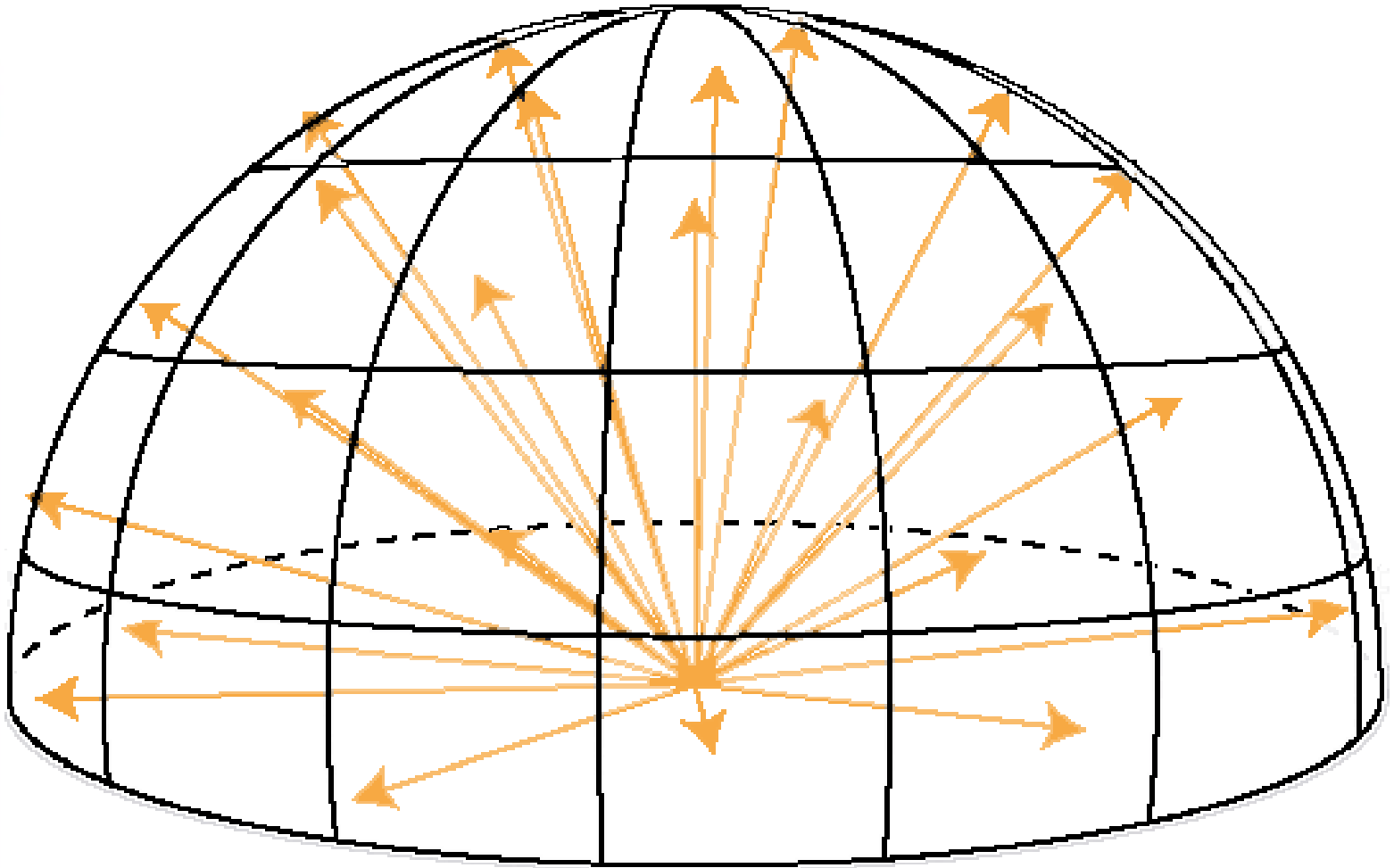
# PRT Coefficient Generation



- G Monte Carlo raytracing method
- r Project, weight, accumulate



# Projecting Visibility Function



# Projecting Visibility Function

```
for all samples {
    stratified2D( u1, u2 ); //stratified random
    numbers
    sampleHemisphere( shadowray, u1, u2, pdf[j] );

    H = dot( shadowray.direction, vertexnormal );

    if ( H > 0 ) { //only use samples in upper
hemisphere
        if ( !occludedByGeometry( shadowray ) ) {
            for ( int k = 0; k < bands; ++k ) {
                RGBCoeff& coeff = coeffentry[k];
                //project onto SH basis
                grayness = H * shcoeffs[j*bands +
project
k];
weight          grayness /= pdf[j];
accumulate     coeff += grayness; //sum up
contribution
            }
        }
    }
}
```



# Projecting Visibility Function

GameDevelopers  
Conference 08

- i 100-400 shadow rays per vertex are sufficient
- a Use your best acceleration structures
- c Vertex normal is accounted for in computation of  $H$ , thus no need for normal in shader
- a Full visibility function is preserved
- f Group objects together in acceleration structure to get inter-object shadowing



# Runtime Reconstruction in Vertex Shader

- u The hard work has been done!
- s Thanks to the orthogonality of SH functions, reconstruction is reduced to a simple dot product of 9 float3's

$$E = \int_s L_i(\mathbf{x}) \cdot BRDF(\mathbf{x}, N) \cdot V(p, \mathbf{x}) dx$$

$$E \cong \sum L_{lm} \cdot P_{lm}$$



# Runtime Reconstruction in Vertex Shader

GameDevelopers  
Conference

08

```
uniform vec3 Li[9];  
attribute vec3 prt0;  
attribute vec3 prt1;  
attribute vec3 prt2;  
  
color = Li[0]*(prt0.xxx) +  
        Li[1]*(prt0.yyy) +  
        Li[2]*(prt0.zzz);  
color += Li[3]*(prt1.xxx) +  
        Li[4]*(prt1.yyy) +  
        Li[5]*(prt1.zzz);  
color += Li[6]*(prt2.xxx) +  
        Li[7]*(prt2.yyy) +  
        Li[8]*(prt2.zzz);
```



CMP  
United Business Media

WWW.GDCONF.COM



# Comparison: With and Without Self-Shadows

GameDevelopers  
Conference 08





# Comparison: With and Without Self-Shadows



# PRT in hyperDrive™



# PRT Compression

Manny Ko  
Naughty Dog



# PRT Compression

- l 3rd order SH occupies 36 bytes per sample
  - u Normals encoded as Lambertian response during pre-processing
  - a World-space => no tangent space needed
  - o Bandwidth is key to GPU performance
- PRT data size is key to pervasive usage in game scenes



# Previous Work in PRT Compression

Sloan introduced CPCA

Break mesh into small clusters

PCA applied to each cluster to obtain a reduced set of bases

- a Fairly effective if the clusters are small, but that will add draw call overhead
- a Hard to apply your vertex cache optimizer



# Our PRT Compression

- i M 1: 9 bytes using scale-bias
  - n M 2: 6 bytes (8; 6,6,6; 6,4,4,4,4)  
M 3: 6 bytes using Lloyd-Max relaxation
  - n M 4: 4 bytes (5; 4,4,4 ;3,3,3,3,3)
- All very simple to implement



# Scalar Quantizer and Scale-Bias

How to convert floats to bytes?

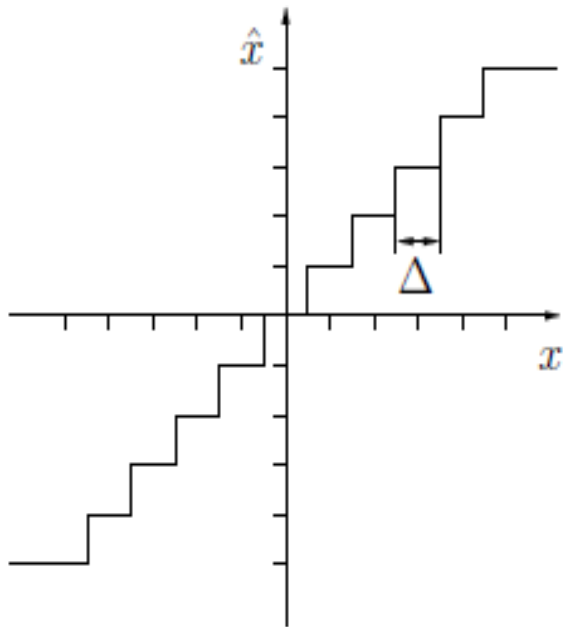
- | Scan each set of coefficients to get range

$$C_i = (P_i - \text{Bias}_i) * \text{Scale}_i?$$

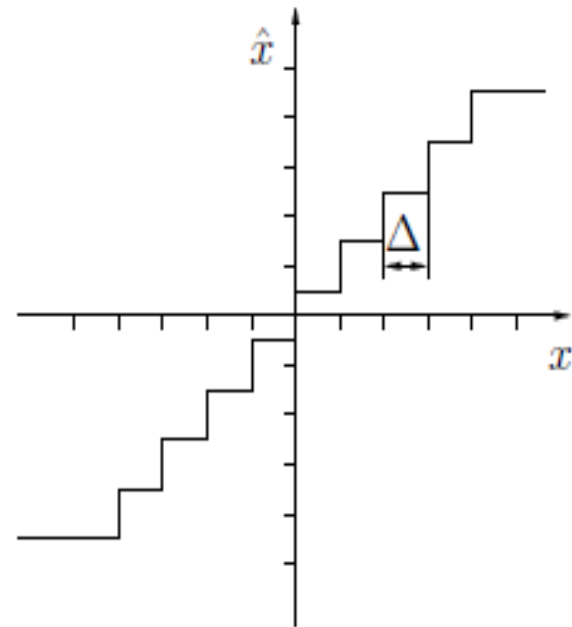
We are dealing with a scalar quantizer



# Mid-rise vs. Mid-thread Quantizer



(a)



(b)



# Mid-rise Quantizer



- e A 2 bit mid-rise quantizer has only 3 decision levels  
Advantage: can construct the DC-level exactly. This can be critical in some cases



# Mid-thread Quantizer

Game Developers  
Conference 08

- i Even number of levels – e.g. 4 levels for a 2 bit code
- e More accurate overall but cannot reconstruct the DC-level
- e For 4 bit codes we improved the PSNR by  $\approx 1$ db



# ScalarQuantizer

```
float half = (qkind ==
PRT::kMidRise) ?
0.5f : 0.f;

for (int i=0; i < nc; i++) {
    float delta = scalars[i] -
bias[i];
    p[i] = delta * scales[i];
    output[i] = floor(p[i] + half);
}
```



# Method 1: 9 Bytes

- s Use Mid-rise quantizer
  - n Remove 4PI factor from coefficients
  - r  $\frac{1}{4}$  the memory and 2X improvement in speed
- No visible quality lost, PSNR > 78db



## Method 2: 48 Bits

- s Method 1 takes us down to 9 bytes. Can we do better?
- s (8,6,6,6,6) (4,4,4,4) – bit fields
- , Choice of bit allocation motivated by
  - Fourier theory => energy compaction
  - shader efficiency
- n Cannot be optimal but try to be close



# Square Norms by Band

## B Energies by band

[0]: 34.8%

[1]: 17.1%

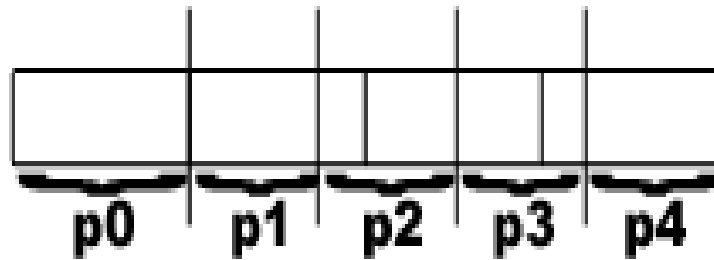
[2]: 15.3%

[3]: 17.1%

85% of energies in 1st 4 bands



1<sup>st</sup> word: 8,6,6,6,6



P0(b31..24),p1(23..18)..



# Vertex Shader: Method 2

- e Coefficients will straddle input registers
- | No bit operation in shaders – use floating point ops to simulate
  - Mul/div by power-of-2 for L/R shifts
  - Fract and trunc() to isolate the 2 pieces. Add for 'or'
  - ) Tries to take advantage of SIMD



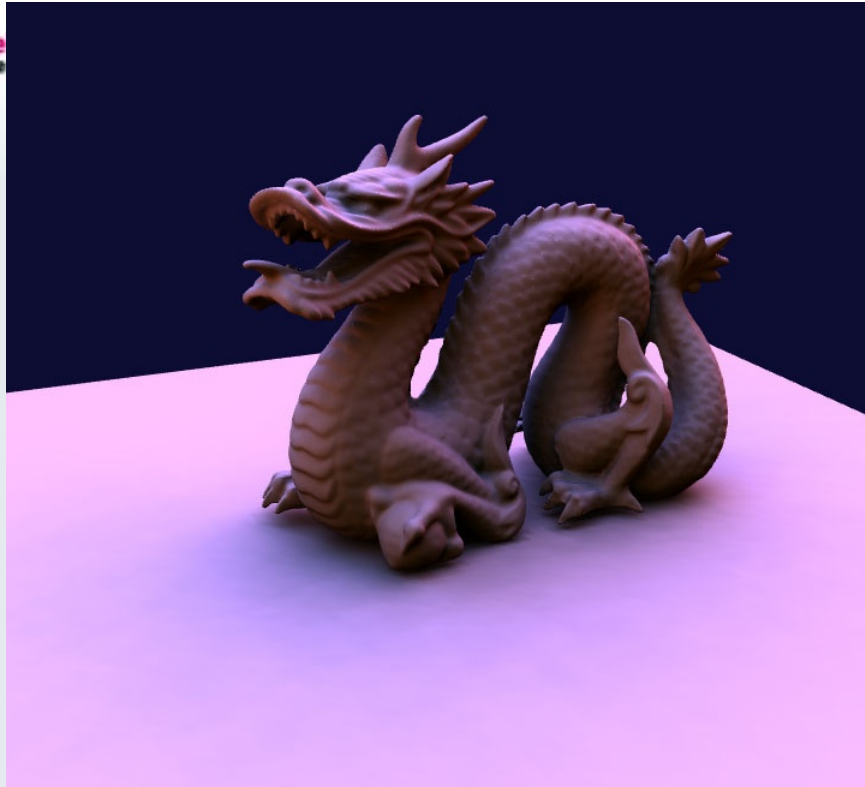
# Vertex Shader

```
0 // rshft(s[0], 1/4, 1/16, 1/64)
/ p14 = prt0 * rshft;
f lhs3.yzw = fract(p14.yzw);
// lshft(0, 16.*4., 4.*16., 64.);
4 lhs3 *= lshft;
4
4 //p[0]:
pp0 = p14.x + bias0;
a //p[1..4]:
a p14.xyz = (lhs3.xyz + floor(p14.yzw)) * scales.xyz + bias.xyz;

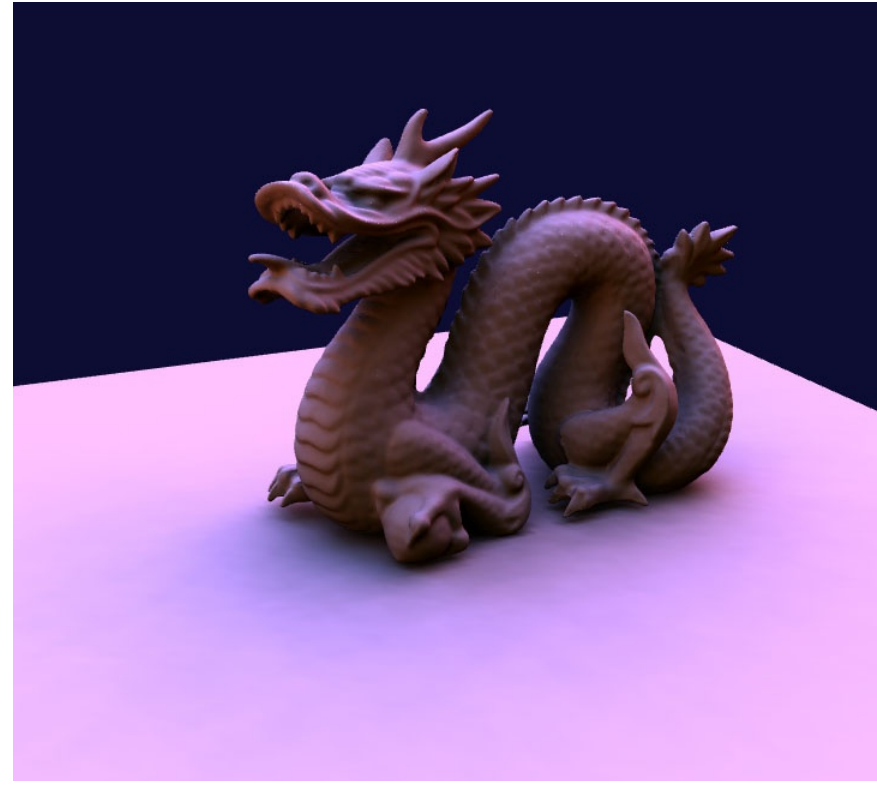
p14.w = lhs3.w * scales.w + bias.w;
```



# Comparison:M1 vs M2



184 fps



243 fps



# Method 2: Demo and Discussion

- n Improved fps by ~17% compared with M1
- ~ Reduced memory usage by 33%. More with alignment restrictions
- s Reduced the number of streams by 1 with room to spare
- e GFLOPS for modern GPU going up much faster than bandwidth
- n Looking for more mathops per byte?



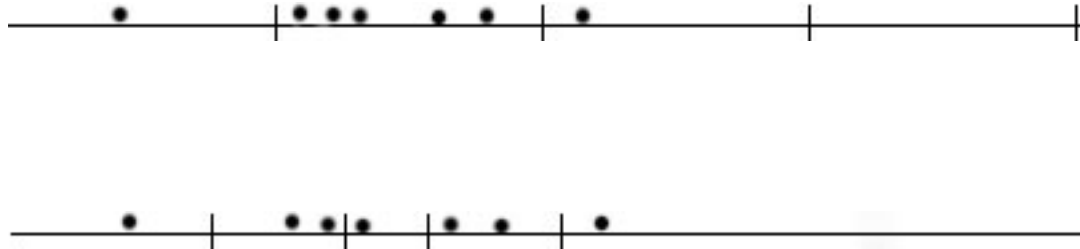
# Uniform vs. Non-Uniform Quantizer

GameDevelopers  
Conference 08

- U A uniform quantizer divides the range evenly
- z Optimal only if all the inputs are equally probable
- a Intuitively in regions of low probability the bin should be wider



# Non-uniform quantizer



How can we design such a set of bins?



# Lloyd-Max Algorithm

- t A version of k-means clustering algorithm
- e Given a fixed number of bins iteratively solves for an optimal set of bins
- i Converges quickly
- y Key is a probability distribution table (*pdf*)  
Applied to all 9 bands separately  
Details in ShaderX6

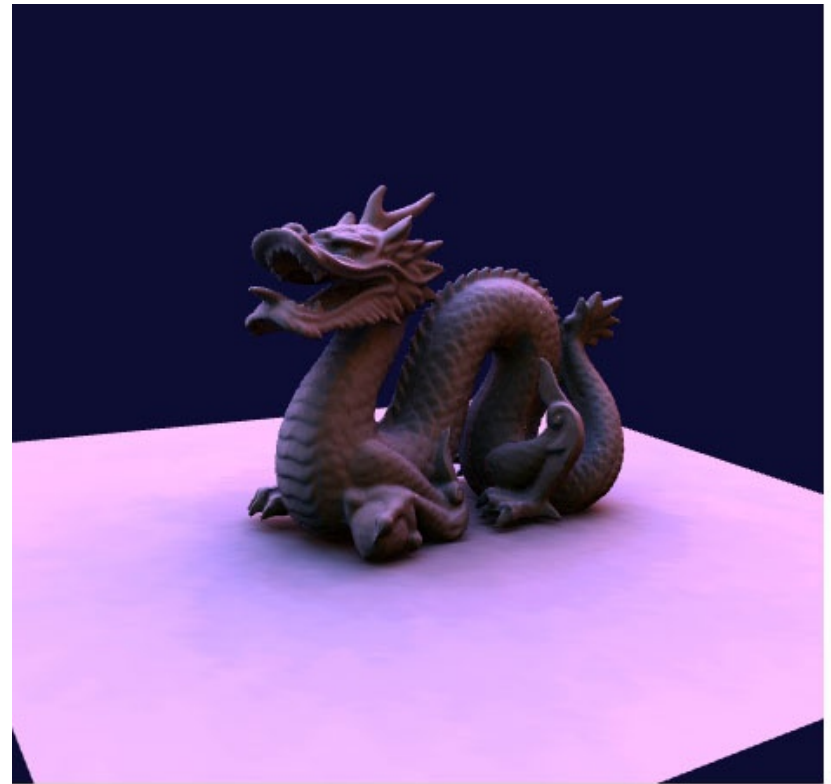
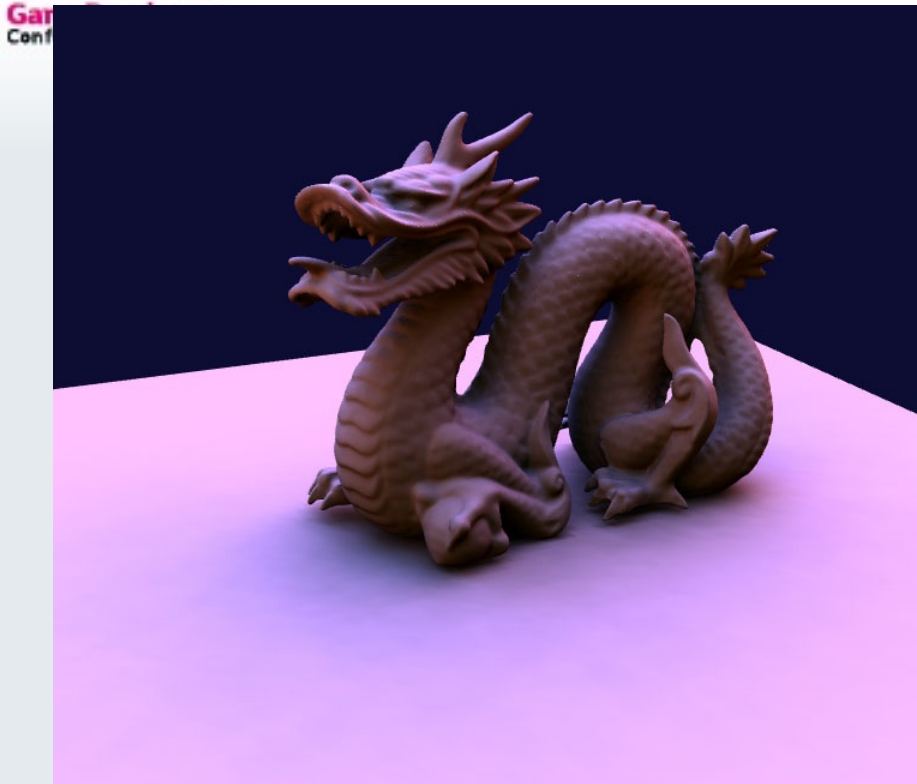


# Lloyd-Max Iteration

- i Given an initial set of bins  $t_k..t_{k+1}$ 
  - Find  $q_k$  such that it is the centroid of  $t_k..t_{k+1}$
  - Find  $t_k$  so that it is the mid-point of  $q_k..q_{k+1}$
- h Repeat for all the bins
- h Check for convergence
- g  $q_k$ s become the reconstruction levels



# Comparison:M2 vs M3





# Quality increase for LM

Improved the PSNR by 1.3 - 1.5db



# Shader for Lloyd-Max

- Almost the same as M2 since we are using the same bit-allocation scheme
- e Decoded bit-fields used to index a table of reconstruction levels `recon [ ]` – the centroids of the bins used in the quantizer
- f Can be constants or a small vertex texture  
Only used for 4 of the 2nd order terms



# Lloyd-Max Summary

- y Fps is about the same or a little slower on older GPUs.  
Same speed on new GPUs
- w Better quality
- w Sensitive to initial condition, see paper for literature references



# Game Scene Demo

Game Developers  
Conference 08





# Limitations

Methods equally good for a map based approach

## Distant Illumination

Good for smaller objects

For large objects either split them or blend multiple *Lis* within shader to avoid seams



# Generalize to Interreflections

- t Instead of sampling visibility, gather indirect illumination
  - i Use a photonmap or irradiance cache or iteratively gather using visibility PRTs
- Use 3x48 or 3x32 bits per sample
- Or pack a 5,6,5 color with one of the 48 bit PRT methods



# Light Sources

GameDevelopers  
Conference 08

- Not limited to environment maps
  - n Pre-integrate any kind of light source – preferably area lights or a large collection of lights, or use probes
- Runtime methods to generate *Lis*
- t Add or subtract lights easily – just add/sub the *Lis*
  - l Rotating lights is cheap



# Normal maps

- Bake the normal variations and fine surface details into the PRT
- ✓ Or use Peter-Pike's [06] method



# Good tool

- , DirectX SDK comes with an offline PRT tool



# Conclusions

- Demonstrate how to implement all key elements of a SH-based PRT shader system
  - Compact and efficient even for older GPUs and no tangent space required
- c Better than simple AO and bent-normal
- l Smaller than bent-normal
- t Groundwork for more GI effects



# More on SH and PRT

GameDevelopers  
Conference 08

- R Peter-Pike's talk Wed. 2:30pm
- k Hao's talk Thu. 4pm
- 4 Yaohua's talk Fri. 2:30pm



# Credits

Matthias Zwicker for his guidance and advising on the PRT research

Ravi for all his help and generosity

Peter-Pike S. for sharing his insights

r Incognito Studio for game assets

Eric H. and Naty H. for figures

Bunkspeed for screenshots

r All our friends



# Contact Info



- Manny Ko – [man961.great@gmail.com](mailto:man961.great@gmail.com)
- i Jerome Ko – [submatrix@gmail.com](mailto:submatrix@gmail.com)



# References

- P [Green 07] “Surface Detail Maps with Soft Self-Shadowing”, *Siggraph 2007 Course*.
- [Ko,Ko 08] “Practical Spherical Harmonics based PRT Methods”, *ShaderX6 2008*.
- X [McTaggart 04] “Half-Life 2 / Valve Source Shading”, *GDC 2004*.
- [Mueller,Haines,Hoffman] *Real Time Rendering 3<sup>rd</sup> Edition*.



# References

- | [Landis02] Production-Ready Global Illumination, *Siggraph 2002 Course*.
- u [Pharr 04] "Ambient occlusion", *GDC 2004*.