# An Improved Shading Cache for Modern GPUs

Pitchaya Sitthi-amorn, Jason Lawrence
University of Virginia

Yang Lei, Pedro V. Sander
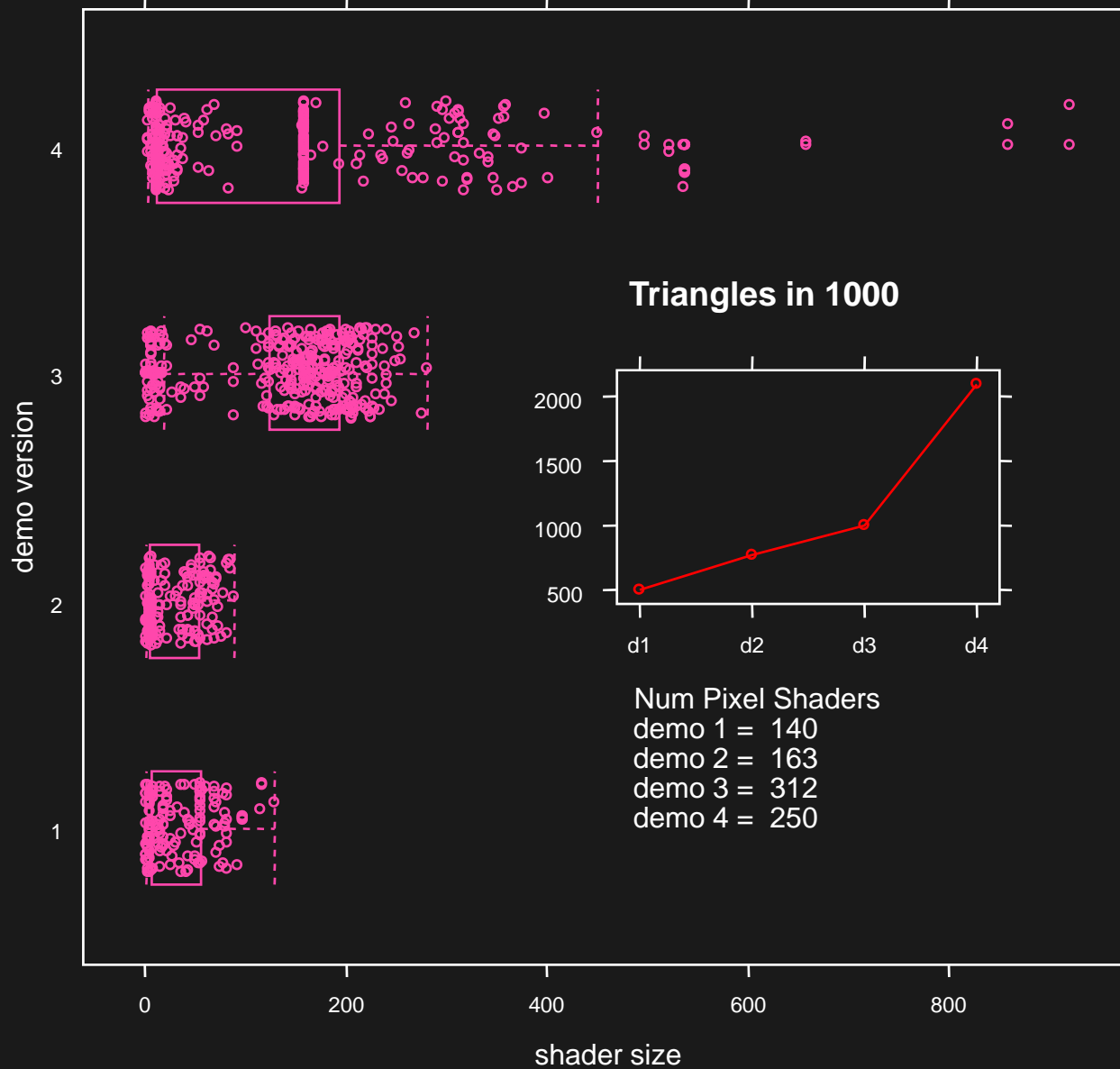Hong Kong University of Science and Technology

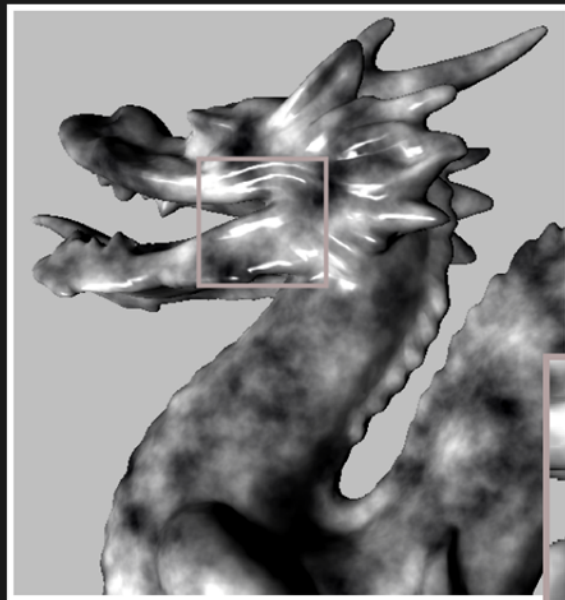Diego Nehab
Microsoft Research

# Motivation

Courtesy of AMD/ATI
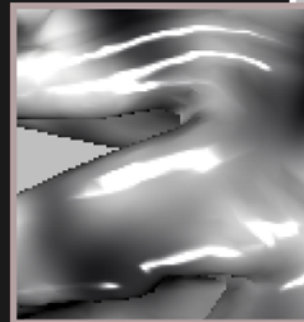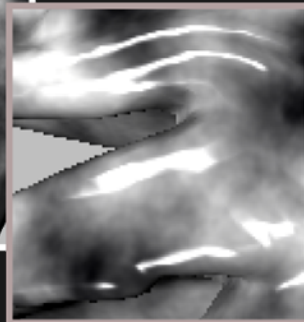
# Shader Complexity of ATI demos

# Related Work: Code Simplification

- Replace subexpressions with constants

- Automatic shader level of detail [Olano et al. 2003]

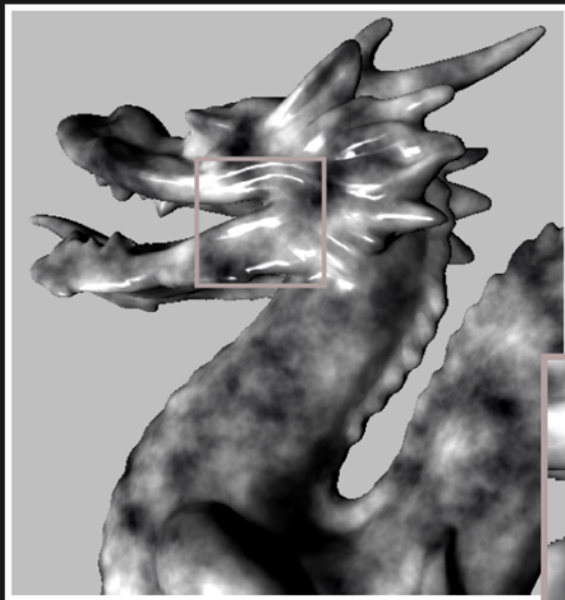- User-configurable automatic simplification  [Pellacini 2005]
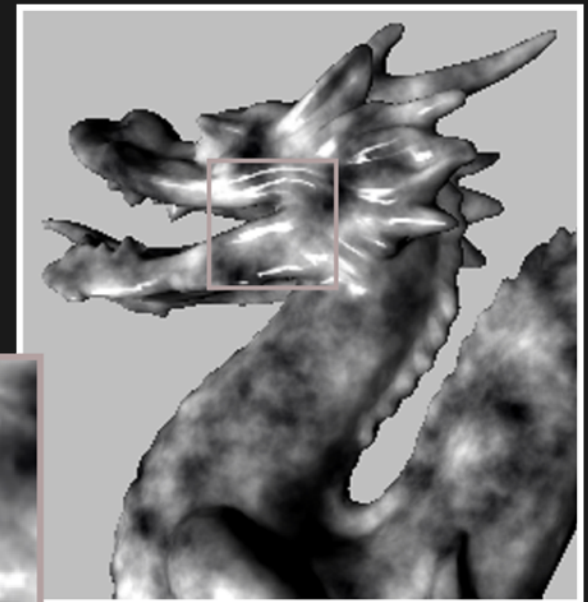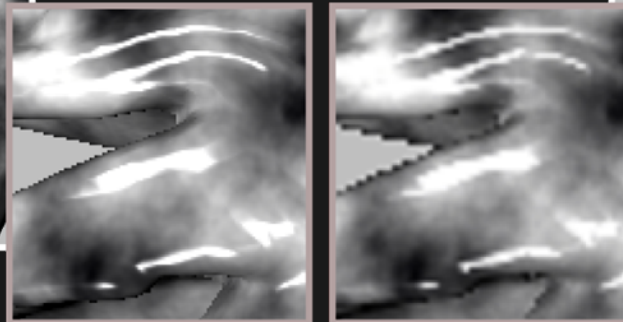


Original

Simplification
(x2.2 faster)

# Related Work: Dynamic Resizing

- Render scene to lo-res off-screen buffer and upsample to target resolution

- *InfiniteReality* system [Montrym et al. 1997]

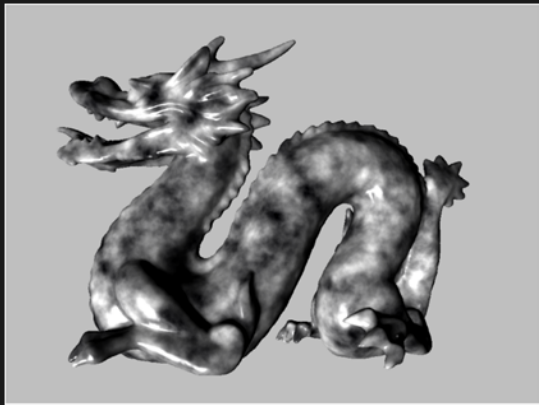- Geometry-Aware resizing [Yang et al. 2008] (concurrent)
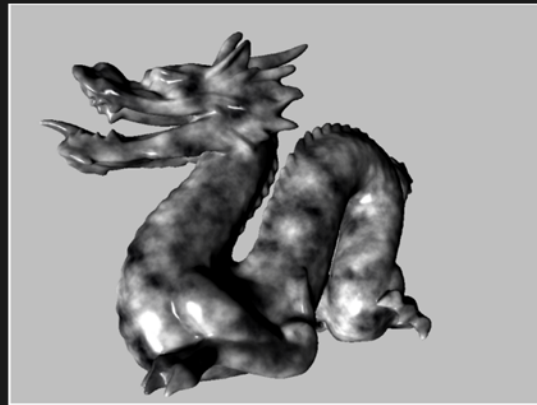


Original

Dynamic Resizing
(x2.7 faster)
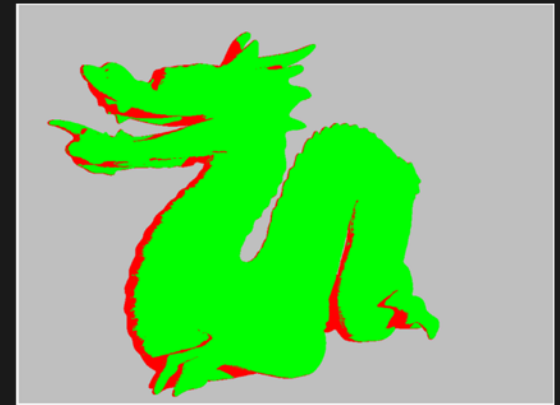
# Related Work: Temporal Reprojection

- Reuse partial shading calculations across consecutive frames

- Reverse reprojection cache [Nehab et al. 2007]

- Pixel-correct shadow maps with temporal reprojection and shadow test confidence [Scherzer et al. 2007]

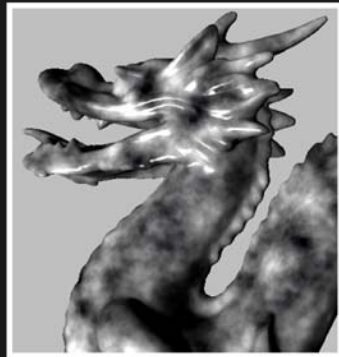- Multi-view architecture [Hasselgren et al. 2006]

frame n-1

frame n

green = mutually visible
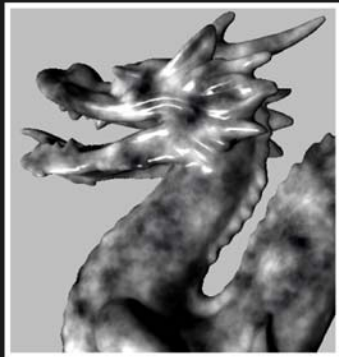red = occluded

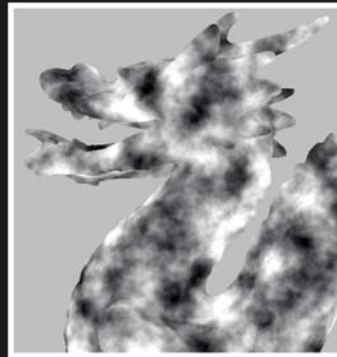# Real Time Shading Cache



Frame n-1

framebuffer

Frame n

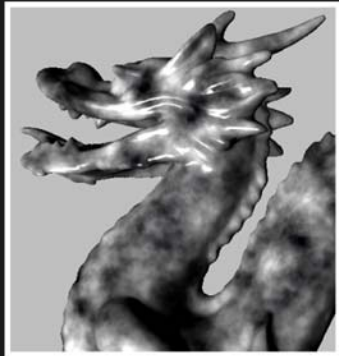# Real Time Shading Cache
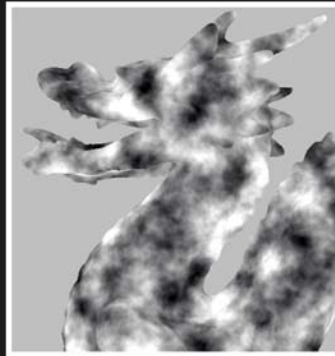


Frame n-1

framebuffer      payload

Frame n

# Real Time Shading Cache



Frame n-1

framebuffer    payload    depth

Frame n

# Real Time Shading Cache



Frame n-1

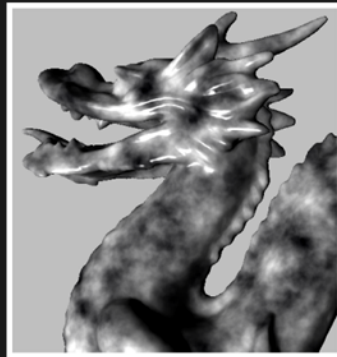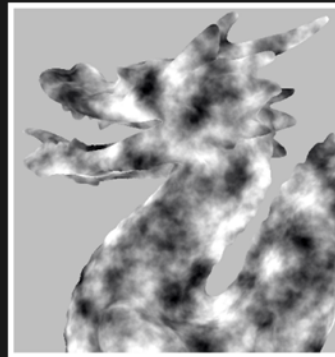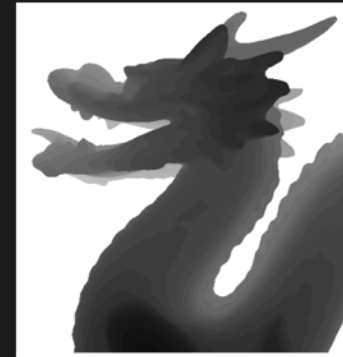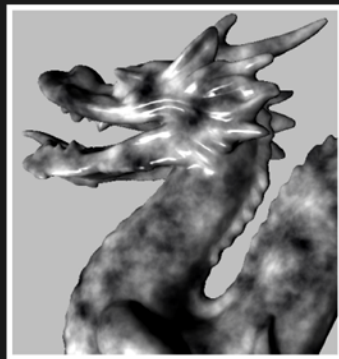framebuffer · payload · depth

Shading Cache

Frame n

# Real Time Shading Cache
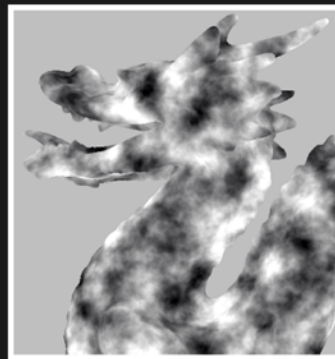


Frame n-1: framebuffer, payload, depth

Shading Cache

Frame n: framebuffer

# Real Time Shading Cache



Frame n-1

framebuffer     payload     depth

Shading Cache

Frame n
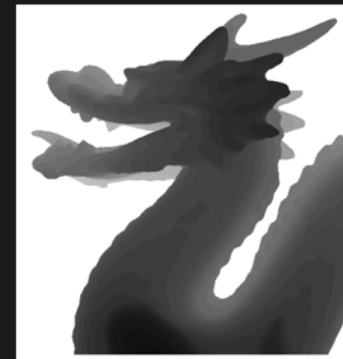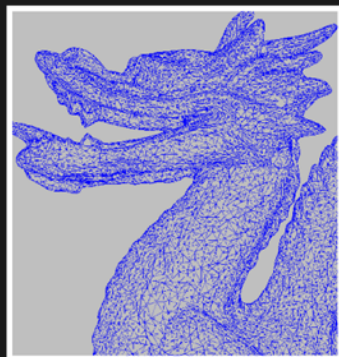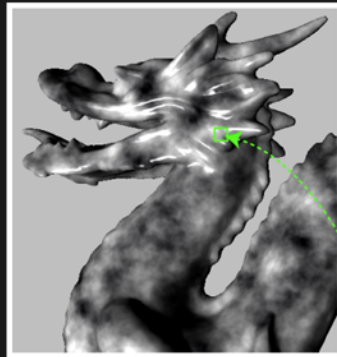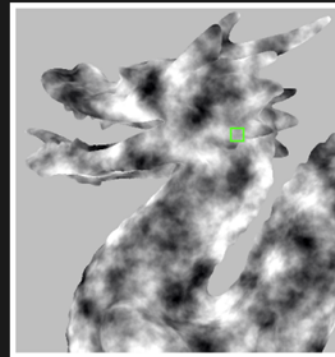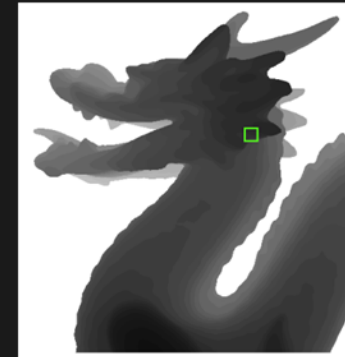
framebuffer

08   11

# Real Time Shading Cache



Frame n-1

framebuffer      payload      depth

Shading Cache

Frame n

framebuffer

08  12

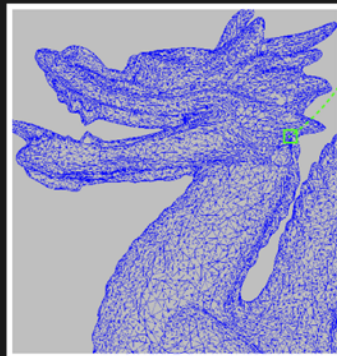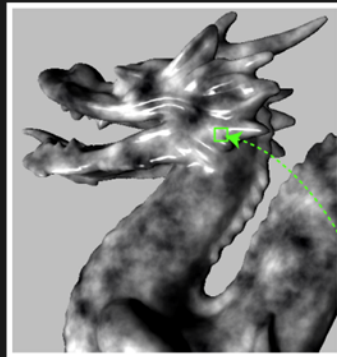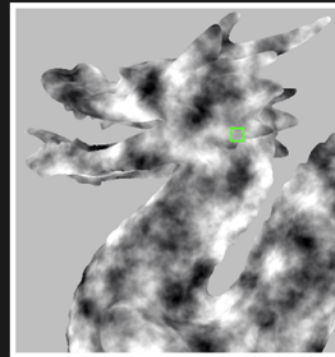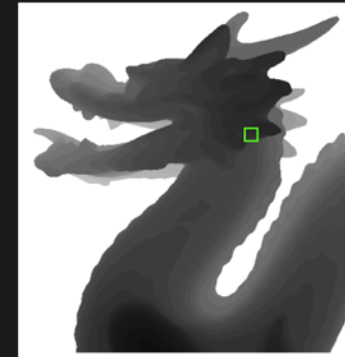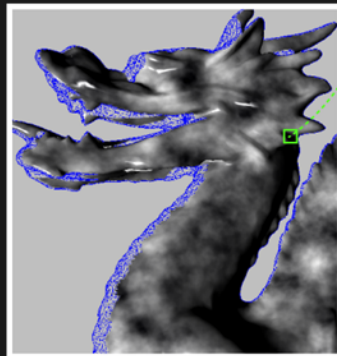# Real Time Shading Cache



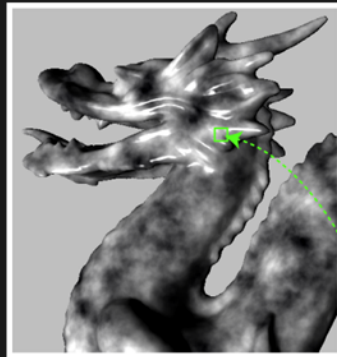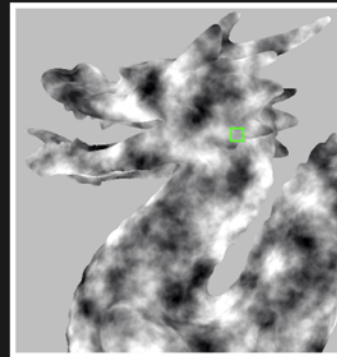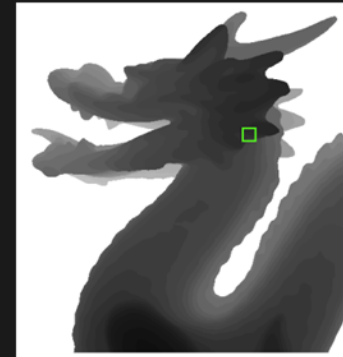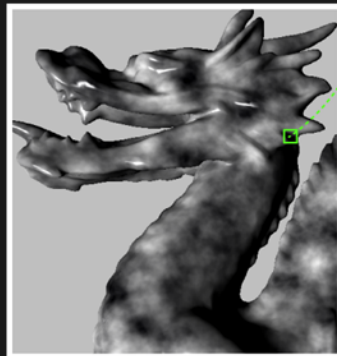Frame n-1 — framebuffer | payload | depth

Shading Cache

Frame n — framebuffer

# Real Time Shading Cache



Frame n-1 — framebuffer, payload, depth

Frame n — framebuffer, payload, depth

Shading Cache

14

# Cache Refresh

- Scene points may remain visible over many frames

- Cached entries will become stale due to changes in shader inputs and from resampling error

- Explicitly refresh cached entries within a user-set refresh period $\Delta n$ by forcing misses within k x k blocks of pixels



$k = 1$  $k = 2$  $k = 4$

$\Delta n = 10$

# Cache Refresh

- Scene points may remain visible over many frames

- Cached entries will become stale due to changes in shader inputs and from resampling error

- Explicitly refresh cached entries within a user-set refresh period $\Delta n$ by forcing misses within k x k blocks of pixels
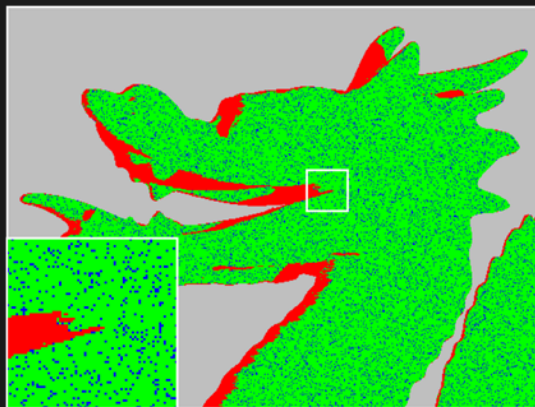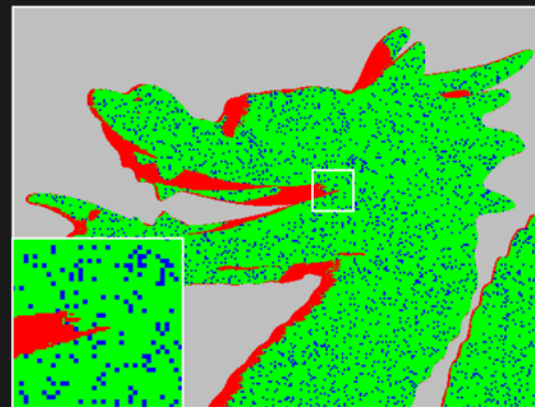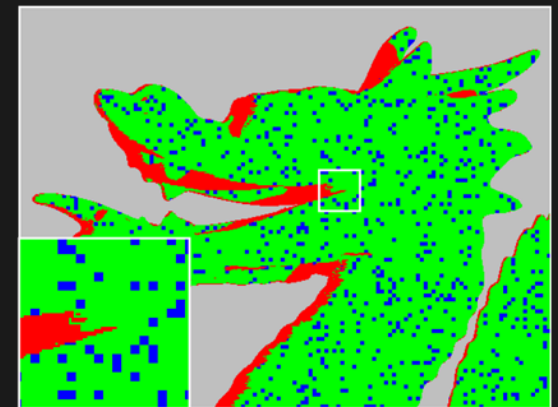


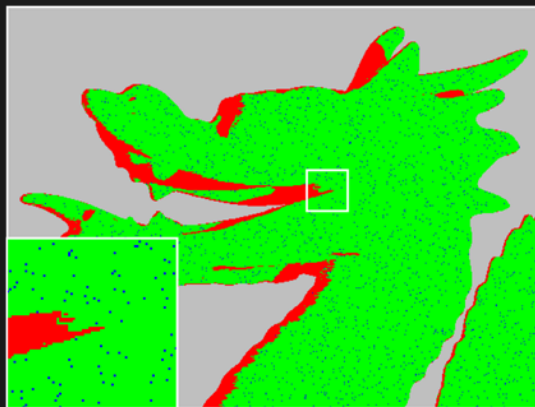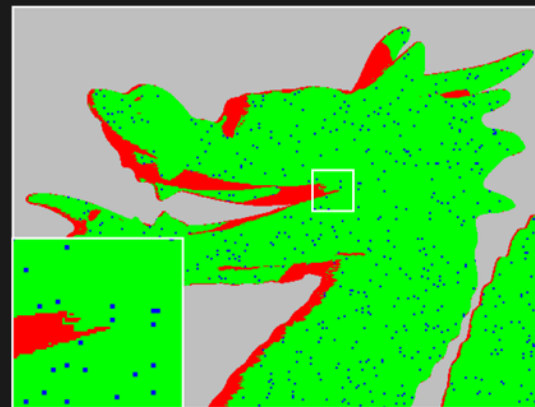$k = 1$        $k = 2$        $k = 4$

$\Delta n = 50$

# 1-Pass Algorithm [Nehab et al. 2007]

# 1-Pass Algorithm [Nehab et al. 2007]



- Branch efficiency of underlying hardware

- Relative cost of processing hit and miss

- Use of multiple render targets (MRTs)

# 2-Pass Algorithm [Nehab et al. 2007]



- Still depends on branch efficiency; however, difference in cost of paths is reduced when hit << miss

- Still requires MRTs

# 3-Pass Algorithm (Our approach)



- Execution paths in the first pass are independent of what is being cached

- Not require MRTs

- Drawback – three rendering passes

# Computation Overlap Problem

# Test Scenes

**Dragon shader**

**Trashcan shader**



procedural noise with
Blinn-Phong specular layer
(75K triangles)

supersampled (25)
environment map
(15K triangles)

# Experiment #1

- Generated versions of the shader that caches every intermediate calculation

- Compute cost of evaluating payload (P)

- Compute cost of evaluating full shader (T)
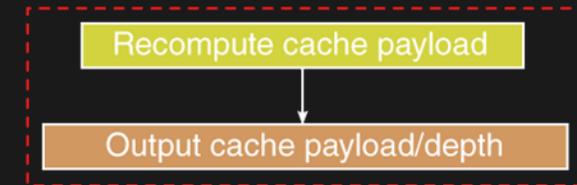
- Fixed refresh period of 32 and 4 x 4 block size

- Compare performance of three different algorithms on NVIDIA Geforce 8800GTX and ATI Radeon 2900TX

NVIDIA GeForce 8800 GTX / Dragon Shader

**P/T increasing**

ATI Radeon 2900 XT / Dragon Shader

# Experiment #1: Trashcan / NVIDIA



NVIDIA GeForce 8800 GTX / Trashcan Shader

Legend:
- 1-Pass
- 2-Pass
- 3-Pass
- Original

Y-axis: Average Render Time (ms)
X-axis: Nodes (sorted by P/T)

**P/T increasing**

ATI Radeon 2900 XT / Trashcan Shader

P/T increasing

Dragon Shader: NVIDIA Geforce 8800GTX

# Experiment #2: Refresh parameters



Dragon Shader: ATI 2900TX

# Conclusion

- Introduced an improved implementation of a shading reprojection cache

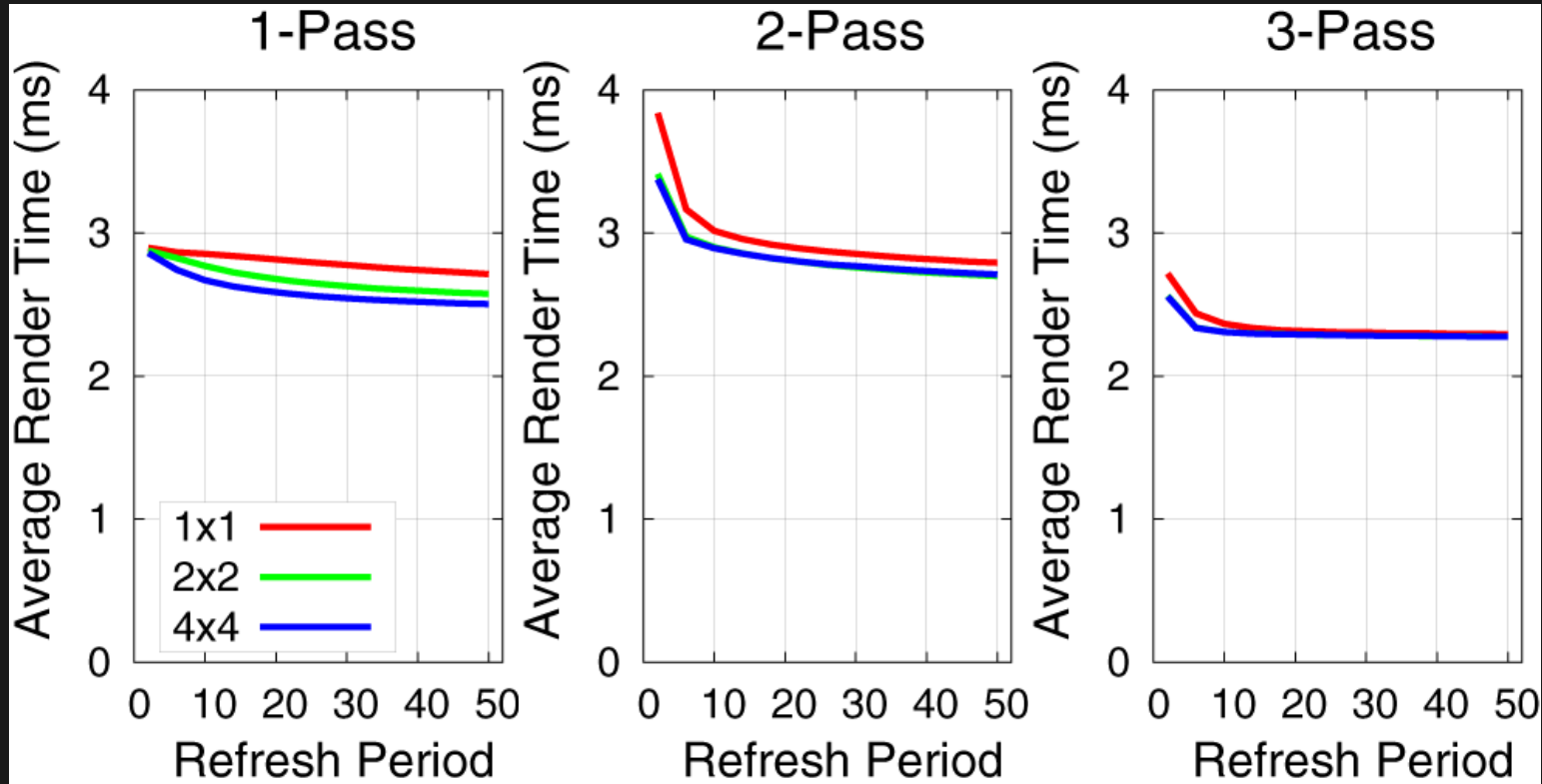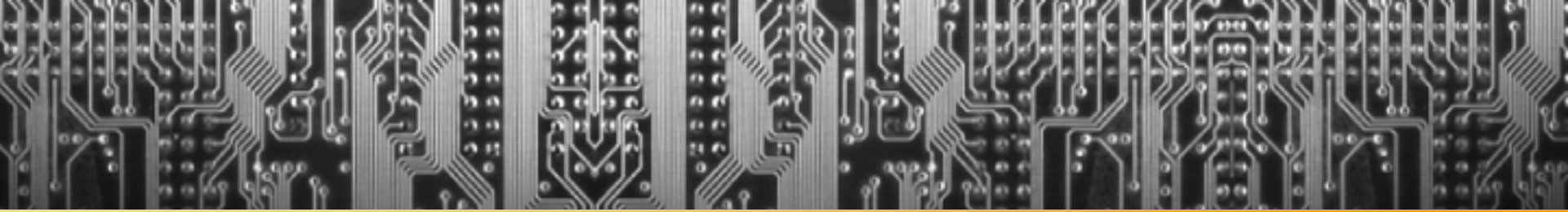- Require single target and limits reliance on efficient branching in hardware

- More consistent performance across a wide range of cache loads on modern NVIDIA and ATI hardware
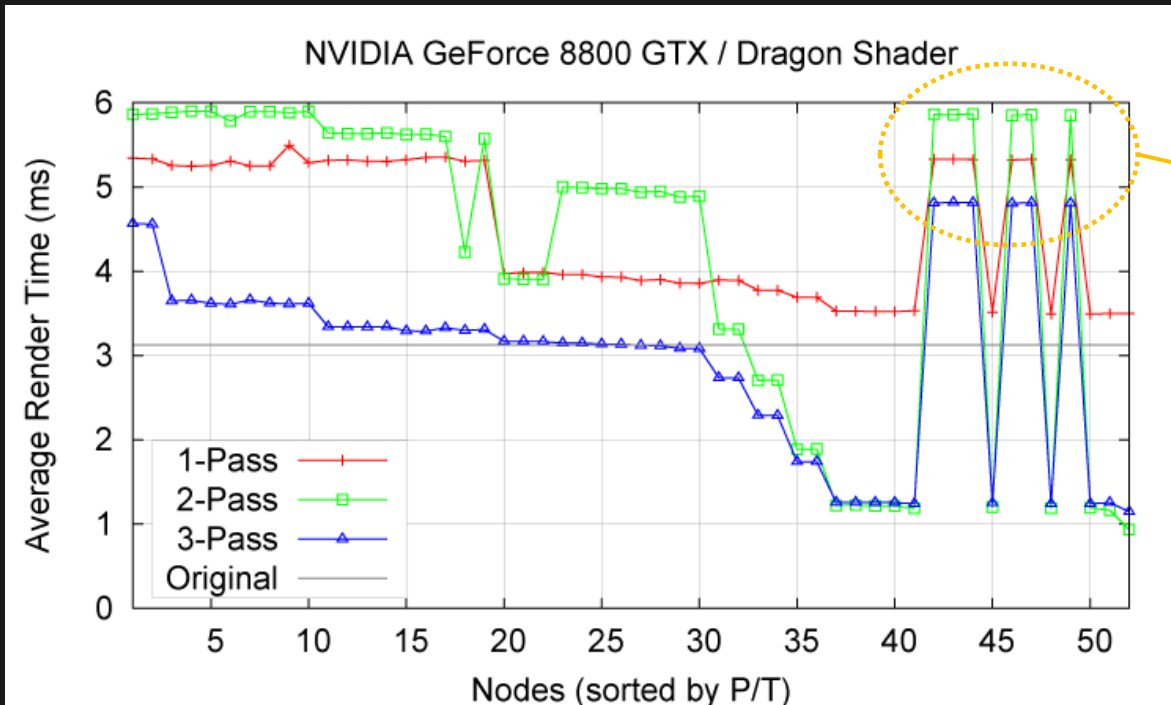
# Future Work

- Explore the possibility of combing existing acceleration techniques

- Automatic cache allocation

- Alternative cache parameterization

# Thank You

# Elaboration on Experiment #1 Results



NVIDIA GeForce 8800 GTX / Dragon Shader

$\alpha$noise()+(1-$\alpha$)noise()

Imagine caching $\alpha$noise() subexpression, noise() would need to be called in both hit and miss paths.