

# GPU-Accelerated High-Quality Hidden Surface Removal

Daniel Wexler, Larry Gritz, Eric Enderton, Jonathan Rice<sup>†</sup>

NVIDIA

---

## Abstract

*High-quality off-line rendering requires many features not natively supported by current commodity graphics hardware: wide smooth filters, high sampling rates, order-independent transparency, spectral opacity, motion blur, depth of field. We present a GPU-based hidden-surface algorithm that implements all these features. The algorithm is Reyes-like but uses regular sampling and multiple passes. Transparency is implemented by depth peeling, made more efficient by opacity thresholding and a new method called z batches. We discuss performance and some design trade-offs. At high spatial sampling rates, our implementation is substantially faster than a CPU-only renderer for typical scenes.*

Categories and Subject Descriptors (according to ACM CCS): I.3.3 [Three-Dimensional Graphics and Realism]: Hidden line/surface removal

---

## 1. Introduction

Rendering for film and other high-end applications typically requires a plethora of hidden surface removal features that have not traditionally been available in graphics hardware: highly supersampled antialiasing, filtered with a high-quality kernel with support wider than one pixel; motion blur and depth of field; order-independent transparency with spectral opacities (separate values for red, green, and blue); and floating-point output of color, opacity, and any other arbitrary values computed by shaders.

Until recently, real-time graphics hardware has taken a simplistic approach to hidden surface removal: a single z-buffer entry per pixel, 8-bit color depth, and an 8-bit alpha channel to represent coverage and opacity. Even now that some of these restrictions have been lifted (notably with the advent of floating-point framebuffers), it has been difficult or impossible to achieve the full array of features above, particularly at the quality levels demanded by high-end applications.

In this paper, we present a high-quality hidden surface removal algorithm that is accelerated by modern commodity graphics hardware. Specifically:

- The algorithm incorporates supersampling, user-selected filters with arbitrarily wide support, depth of field, and multiple output channels.
- The algorithm supports transparency. It uses the depth peeling technique, enhanced to allow multi-channel opacity and opacity thresholding. We also present optimizations that, for typical scenes, allow transparency performance to scale linearly rather than as  $O(N^2)$ .
- The algorithm uses grid occlusion culling to avoid additional shading.
- The algorithm produces comparable quality and superior performance to the hidden surface removal algorithms used in CPU-only software renderers, for many typical scenes (see Figure 1).



**Figure 1:** A motion blurred image from our GPU-accelerated renderer. Courtesy of Tweak Films.

- We explore a variety of engineering trade-offs in implementing this algorithm on modern GPUs, and discuss its performance on a variety of scenes.

This work does not attempt to create real-time rendering capabilities (e.g., 30 fps). Rather, our goal is to accelerate the rendering of images at the highest quality possible, suitable for film or broadcast, that now take minutes or hours. We have already achieved significant speed improvements by using graphics hardware, and we expect those improvements to become even greater over time as graphics hardware becomes more capable and its speed improvements continue to outpace those of CPUs.

### 1.1. Rationale: Programmable Graphics Hardware

Previous attempts at high-quality rendering with graphics hardware have generally involved costly custom hardware

---

<sup>†</sup> {dwexler|lgritz|eenderton|jrice}@nvidia.com

whose performance is often soon overtaken by improvements in general purpose CPUs.

In contrast, this work is concerned strictly with commodity graphics hardware. We find this approach promising because the economy of scale present with modern commodity hardware makes it very inexpensive, and because graphics hardware increases in capability at a much faster rate than CPUs are improving. Kirk [Kir02] reports that for several years graphics hardware doubled in speed every six to twelve months versus every eighteen months for CPUs.

In addition to traditional geometric transformation and rasterization, modern graphics processing units (hereafter referred to as “GPUs”) such as the ATI Radeon 9800 and NVIDIA GeForce FX feature (1) significant programmability via vertex and fragment shaders, rapidly converging on the capabilities of stream processors; and (2) floating-point precision through most or all of the pipeline [ATI03, NVI03]. These two facilities are the key enablers of the work described here.

## 1.2. Related Work

Prior work can be categorized according to the hardware it requires. Software-only systems for high-quality rendering, particularly the Reyes architecture [CCC87, AG99] and its commercial implementations, form both the jumping-off point for our approach and the baseline against which we measure our results. Systems using specialized hardware for high-quality rendering have been proposed or built, but are outside the scope of this paper. This leaves GPU and GPU-assisted methods.

Current GPUs support *multisampling*, which computes visibility for multiple samples per pixel but reuses a single color from the center of the pixel [MH02]. This improves real-time graphics but is not high-quality: current hardware limits the filter shape to a box with no overlap between adjacent pixels, full floating-point formats are not yet supported, and only low, fixed numbers of samples per pixel are supported (currently 4 to 8, versus the dozens typically used for film rendering). The proposed Talisman architecture [TK96] mitigates the memory expense of multiple samples by rendering one screen region at a time, and also supports multisampling, though at a fixed  $4 \times 4$  samples per pixel.

The *accumulation buffer* [HA90] supports antialiasing, motion blur, and depth of field effects by accumulating weighted sums of several passes, where each pass is a complete re-rendering of the scene at a different subpixel offset, time, and lens position, respectively. Direct hardware support has not yet included full 32-bit floating point precision, but programmable GPUs can use fragment shaders to apply this technique at floating point precision, although still at the expense of sending geometry to the GPU once per spatial filter sample. Interleaved sampling [KH01] can reduce the artifacts from regularly-sampled accumulation buffers.

Scenes with partially transparent surfaces, possibly mixed with opaque surfaces, are challenging. Methods that require sorting surfaces in depth are problematic because intersecting surfaces may have to be split to resolve depth ordering. The *depth peeling* method, described in general by Mammen [Mam89] and for GPUs by Everitt [Eve01], solves *order-independent transparency* on the GPU. The method is summarized later in this paper. Its principal drawback is that a scene with maximum depth complexity  $D$  must be sent to the GPU  $D$  times. For some scenes, such as particle systems for smoke, this can result in  $O(N^2)$  rendering time, where  $N$  is the number of primitives. Mammen suggests an *occlusion culling* optimization (though not by that name). In his algorithm, once an object no longer affects the current pass, it is dropped for subsequent passes. Kelley et al. [KGP\*94] designed hardware with four  $z$ -buffers, able to depth peel four layers at once. They use screen regions so that only regions with high  $D$  require

high numbers of passes. The hardware *R-Buffer* [Wit01] proposes to recirculate transparent pixels rather than transparent surfaces; this avoids sending and rasterizing the geometry multiple times, but can require vast memory for deep scenes.

## 2. Architectural Overview

The context for this work is a Reyes-like architecture [CCC87, AG99]. High-order surface primitives (such as NURBS, subdivision surfaces, and polygons) are recursively split until small enough to shade all at once. These patches are discretized into *grids* of pixel-sized quadrilateral micropolygons. Grids typically contain on the order of 100 micropolygons, and are limited to a user-selected size. The grids are possibly displaced, then assigned color, opacity, and optionally other data values at each grid vertex by execution of user-programmable shaders. Finally, the shaded grids are undergo hidden surface removal (“hiding”) to form a 2D image.

It is this final stage—hidden-surface removal of the shaded grids of approximately pixel-sized quadrilaterals—that is the concern of this paper. Whether performed on the GPU or CPU, the earlier stages—handling high-order geometry and shading the micropolygon vertices—are largely orthogonal to the methods used for hidden surface removal. Although beyond the scope of this paper, our shading system implements advanced techniques including global illumination and ray tracing.

Like many Reyes-style renderers, we divide image space into rectangular subimages called *buckets* [CCC87] in order to reduce the working set (of both geometry and pixels) so that scenes of massive complexity can be handled using reasonable amounts of memory. For each bucket, the grids overlapping that bucket are rasterized using OpenGL, as described in Section 3. To handle partially-transparent grids, we use depth peeling but extend it in several important ways to allow multi-channel opacity and opacity thresholding. We also introduce a batching scheme that reduces the computational complexity of depth peeling for typical cases. The transparency algorithms are discussed in detail in Section 4.

Motion blur and depth of field are achieved using an accumulation-buffer-like technique involving multiple rendering passes through the geometry for the bucket. This is described in Section 5.

In order to achieve sufficiently high supersampling to ameliorate visible aliasing, buckets are rendered for all of these passes at substantially oversampled resolution, then filtered and downsampled to form the final image tiles. The filtering and downsampling is performed entirely on the graphics card using fragment programs, as described in Section 6.

The handling of transparency, motion blur, depth of field, and arbitrary output channels can lead to large numbers of rendering passes for each bucket (though much of that work must be performed only for buckets that contain transparent or moving geometry). Specifically,

$$\text{passes} = P_{\text{output}} \times P_{\text{motion}} \times P_{\text{transparent}}$$

where  $P_{\text{output}}$  is the number of output images;  $P_{\text{motion}}$  is the number of motion blur or depth of field samples; and  $P_{\text{transparent}}$  is the number of passes necessary to resolve transparency.

GPUs are well optimized to rasterize huge amounts of geometry very rapidly. Large numbers of passes render quite quickly, generally much faster than using CPU-only algorithms. In short, brute force usually wins. We will return to discussion of performance characteristics and other results in Section 7.

## 3. Hiding Opaque Surfaces

This section describes our basic hidden-surface algorithm for opaque geometry. For simplicity, we describe the process for

a single bucket in isolation; extension to multiple buckets is straightforward. The steps are repeated for each output image<sup>†</sup>.

The algorithm has the overall splitting-and-dicing form of a Reyes-style algorithm, but with a new grid occlusion test towards the end:

```
for each output image:
  for every object from front to back:
    if object's bounding box passes the occlusion test:
      if object is too big to dice:
        split object and re-insert split sub-objects
      else:
        dice object into grid of final-pixel-sized quads
        if diced grid passes the occlusion test:
          shade the grid
          render the grid
```

The grid rendering and both occlusion culling steps are executed on the graphics hardware.

### 3.1. Grid Rendering

To render diced grids with high quality on the GPU, we use *regular supersampling* — in other words, we render large, and minify later. We render the grids using standard OpenGL primitives into an image buffer which is larger than the final pixel resolution by a user-selected factor, typically 4-16x in each dimension. For a typical bucket size of  $32 \times 32$  pixels, the supersampled buffer fits easily in GPU memory. It is a happy synergy that buckets, primarily designed to limit in-core scene complexity, also help to limit buffer size. Rendering grids large also helps tune the geometry for the GPU. The projected screen size of a grid micropolygon is usually on the order of a single pixel in the final image. However, current GPUs are designed to be most efficient at rendering polygons that cover tens or hundreds of pixels; smaller polygons tend to reduce the parallel computation efficiency of the hardware. Supersampling helps to move our grid polygons into this “sweet spot” of GPU performance.

Shading data can be computed at the resolution of the grid—the *dicing rate*—or at another, higher resolution. If the shading and dicing rates match, then the shaded colors are passed down to the GPU as per-vertex data. A full-precision register interpolates the colors across polygons, resulting in smooth-shaded grids.

If the shading rate differs from the dicing rate, so that shading data is no longer per-vertex, then colors are passed to the GPU as a floating-point texture instead. Sending the texture to the GPU and then sampling it for each fragment is slower than passing a color per vertex. In our implementation, the penalty was 10–20%, and quality suffered because the hardware we used could not efficiently filter floating-point texture lookups.

### 3.2. Occlusion Culling

In modern Reyes-style algorithms, before an object is split, diced, or shaded, its bounding box is tested to see whether it is completely occluded by objects already rendered [AG99]. If so, it is culled. To maximize culling, objects are processed in roughly front-to-back order, using a heap data structure ordered by the near  $z$  values of the objects’ camera space bounding boxes. In a high-quality renderer, shading tends to be more expensive than occlusion testing, since it may include texture lookups, shadows, or even global illumination. So occlusion culling before shading can often produce great performance benefits.

CPU-based renderers typically maintain a hierarchical  $z$ -buffer for occlusion testing [GK93]. By contrast, our renderer uses the GPU hardware to occlusion test against the full  $z$ -buffer, via the OpenGL *occlusion query* operations [SA03]. Occlusion query returns the number of fragments that passed the  $z$ -buffer test between the query start and end calls. In other words, it returns the number of visible fragments. We turn off writes to the framebuffer, begin an occlusion-query, render the bounding box, then end the query. If the query reports zero visible fragments, the object is culled.

The bounding box occlusion cull before objects are diced is similar to CPU-based algorithms. But later, after the object is diced (but before it is shaded), we occlusion cull again using the actual grid geometry, an exact test for whether rendering the grid would change any pixels. This test would be quite expensive on the CPU, but GPUs rasterize quickly. Grid culling is especially effective near silhouette edges, where grids behind the silhouette very often have a bounding box that spills out past the silhouette. Grid culling can also help when an object’s bounding box pokes out from behind other geometry, but the actual object does not (e.g., two concentric spheres with slightly different radii). Grid culling reduces shading by 10–20% in our test scenes and provides a significant performance boost. In scenes dominated by shading time such as ambient occlusion, ray tracing, and complex user-defined shaders, the performance boost is even more pronounced.

The PC bus architectures in current use allow much greater data bandwidth to the GPU than back from it. Occlusion-query fits this mold well, since it returns only a single integer. However, the GPU must finish rendering the primitives before the occlusion-query result is available. Thus the frequent queries in our hider algorithm tend to cause the GPU to run at reduced efficiency, although some latency can be hidden by careful ordering of operations.

## 4. Transparency

As mentioned earlier, we use depth peeling to render grids that may be transparent and may overlap in depth and screen space. Standard depth peeling [Eve01] is a multipass method that renders the nearest transparent surface at each pixel in the first pass, the second nearest surface in the second pass, etc. Each pass is composited into an RGBAZ “layers so far” buffer that accumulates the transparent layers already rendered. Each pass renders the whole scene using ordinary  $z$ -buffering, with an additional test that only accepts fragments that are behind the corresponding  $z$  in the layers-so-far buffer. An occlusion query test on each pass indicates whether that layer was empty; if it was, the algorithm halts.

We considered an alternative to depth peeling, an algorithm to “bust and sort” overlapping grids into individual micropolygons, then render them all in front-to-back order. This approach was abandoned because current GPUs do not support floating-point alpha blending or allow framebuffer reads from within fragment programs. This might be worth revisiting if future GPUs add these capabilities.

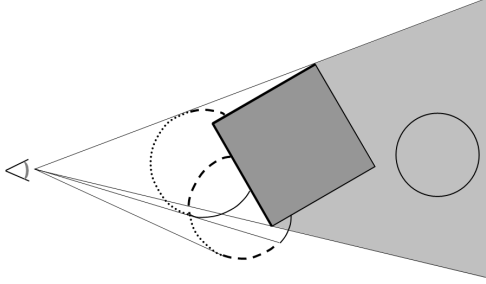
Our depth peeling algorithm (Listing 1) extends that of Everitt in several ways. First, like Mammen [Mam89], we process all opaque surfaces beforehand into a separate RGBAZ buffer, leaving only the transparent grids for depth peeling. Second, we perform opacity thresholding between batches of transparent surfaces sorted in depth, resulting in greatly improved performance for typical cases. Third, we handle spectral opacity, using three passes when needed.

### 4.1. Opaque Preprocessing

Because typical scenes are mostly opaque, rendering opaque surfaces first drastically cuts the number of required depth peeling passes. The remaining transparent grids are occlusion culled against the opaque  $z$ -buffer, further reducing depth peeling effort. We use the opaque  $z$ -buffer texture as an additional

<sup>†</sup> Users’ shaders can optionally output other variables besides color. Each is rendered into a separate *output image*. Examples include surface normals or separate specular and diffuse images.

depth comparison during depth peeling, so transparent micropolygons occluded by opaque surfaces do not add additional depth peeling passes.



**Figure 2:** Depth-peeling of three transparent spheres and an opaque cube, after one layer has been peeled. Dotted lines indicate the layers-so-far buffer (nearest layer) while dashed lines indicate the layer being computed (second-nearest layer). The thick line on the surface of the cube indicates the opaque-surfaces buffer, which culls away all surfaces behind it.

Figure 2 illustrates the algorithm. The opaque preprocess renders all opaque surfaces into the opaque  $z$ -buffer. From here on, we render only transparent surfaces. Grids occluded by the opaque  $z$ -buffer are culled. Pass 1 computes RGBAZ of the nearest transparent surfaces; this initializes the layers-so-far buffer. Pass  $p$ ,  $p > 1$ , computes RGBAZ of the  $p$ -th nearest transparent surfaces. The fragment program for Pass  $p$  rejects fragments unless they are both behind the layer-so-far buffer's  $z$  and in front of the opaque buffer's  $z$ ; the  $z$ -buffer test for Pass  $p$  selects the nearest of the accepted fragments. We halt when the occlusion query test for Pass  $p$  reports that no fragments were accepted. Otherwise, we composite the RGBA of Pass  $p$  under the layers-so-far buffer, and replace the  $z$  of the layers-so-far buffer with the  $z$  of Pass  $p$ .

#### 4.2. Z-Batches

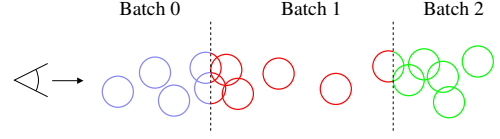
Depth peeling has  $O(N^2)$  worst-case performance for  $N$  grids. If we could depth peel the grids in batches of  $B$  grids each, the worst case would be only  $O((N/B)B^2) = O(BN) = O(N)$ . The problem is that the grids from separate batches may overlap in depth, so we cannot simply depth peel the batches independently and composite the results. We solve this problem by restricting each batch to a specific range of  $z$  values, partitioning the  $z$  axis so that standard compositing yields correct results.

We do this with constant- $z$  clipping planes. Recall that we have sorted our grids by the  $z$  value of the front plane of the camera-space bounding box. We break the list of transparent grids into  $z$ -batches of  $B$  consecutive grids each. The  $z_{\min}$  of each  $z$ -batch is the  $z_{\min}$  of its first primitive. While depth peeling a  $z$ -batch, we clip away fragments nearer than the batch's  $z_{\min}$  or farther than (or equal to) the next batch's  $z_{\min}$ ; see Figure 3. Grids that cross a  $z_{\min}$  boundary are rendered with both batches. Each  $z$ -batch now produces the correct image for a non-overlapping range of  $z$  values, and simple compositing of these images now works.

Grids that cross multiple  $z_{\min}$  planes appear in multiple  $z$ -batches. In the worst case, grids overlap so extensively in depth that  $z$ -batch size effectively approaches  $N$ , and we still do  $O(N^2)$  work. More typically, batches are bigger than  $B$  but still much smaller than  $N$ , and the speed-up is enormous.

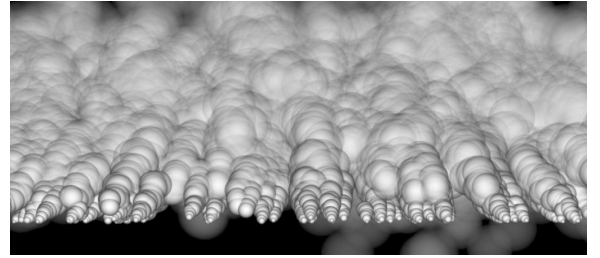
#### 4.3. Opacity Thresholding

Due to the nature of the Porter-Duff *over* operation [PD84] used to accumulate transparent layers, the opacity value will approach but never equal full opacity. Opacity thresholding



**Figure 3:** Clipping planes in  $z$  let us depth-peel each  $z$ -batch independently. Here we see three  $z$ -batches, with  $B = 5$ . Assume each sphere is one grid. Note that grids that “trail” into a following batch are rendered in both batches, with  $z$ -clipping at the batch boundary.

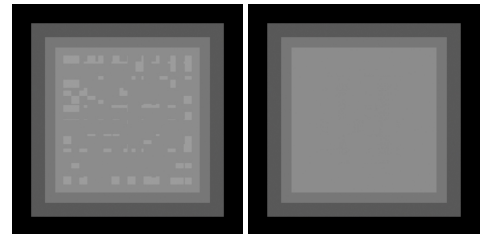
[Lev90] approximates opacity values above a user-specified threshold as fully opaque, thus reducing the effective depth complexity. This meshes nicely with  $z$ -batches. After each  $z$ -batch is rendered, we identify these pseudo-opaque pixels and merge their  $z$  values into the opaque buffer. Before grids are added to a  $z$ -batch, they are occlusion culled against this opaque buffer.



**Figure 4:** A particle system with 16,000 transparent overlapping spheres. Courtesy of Tweak Films.

Using  $z$ -batches and opacity thresholding can improve performance by orders of magnitude. Figure 4 rendered in just under 4 hours without using either technique; in 13 minutes with  $z$ -batching ( $B = 26$ ) but no opacity thresholding; and in 3 min. 38 sec. with  $B = 26$  and an opacity threshold of 0.95. The stochastic CPU renderer discussed in Section 7 renders the scene in 5 min. 15 sec. with the same opacity threshold.

Thresholding results are inexact because the threshold test is performed after an *entire*  $z$ -batch is processed and not as individual fragments are processed. The first grid to push a pixel over the threshold may be in the middle of a batch, so the pixel may still accumulate a few more layers. This artifact is easy to see when we lower the opacity threshold to an unreasonably low value such as 0.2, especially since batching can change at bucket boundaries (see Figure 5). However, it is nearly invisible at ordinary opacity thresholds (e.g., 0.95).



**Figure 5:** Exaggerated artifacts of opacity thresholding by batches. A stack of camera-facing planes, each with opacity 0.1, with the opacity threshold set to the unreasonably low value of 0.2. Image on right shows the ideal solution.

If future GPUs were to extend the occlusion query operation to return the pixel bounding box and the  $z$  range of the rendered fragments that passed the occlusion test, we could optimize the thresholding pass by running it only on the bounding rectangle from the previous depth peeling pass. This would allow us to run the test after each primitive instead of per  $z$ -batch, which would reduce error and improve culling efficiency. Occlusion query is the only reduction operator currently available in GPU hardware; enhancements could be useful for many general computations.

#### 4.4. Spectral Opacity

During geometry processing, we detect whether any primitive contains spectral opacity, as opposed to monochromatic opacity. Hardware supporting *multiple draw buffers* [ATI03] can render spectral opacity in a single pass. Otherwise, three passes of the whole transparency algorithm are required, one each for red, green, and blue.

---

**Listing 1** Pseudo-code for the transparency rendering loop including depth peeling,  $z$ -batch management and opacity thresholding.

---

```

for every object from front to back:
    if object's bounding box passes the occlusion test:
        if object is splittable:
            split object and re-insert split sub-objects
        else:
            dice object into grid
            if diced grid passes the occlusion test:
                shade grid
            if grid is transparent:
                append grid to current  $z$ -batch
                if current  $z$ -batch size exceeds threshold:
                    save opaque RGBAZ
                    while not finished rendering all transparent layers:
                        for every transparent grid in  $z$ -batch:
                            render transparent grid
                        composite transparent layer under layers so far
                        store accumulated transparent RGBA
                        restore opaque RGBAZ
            else:
                render opaque grid

if transparent grids were rendered:
    composite accumulated transparent layers over opaque layer

```

---

#### 5. Motion Blur and Depth of Field

Motion blur and depth of field are computed by rendering multiple passes into a supersampled accumulation buffer stored on the GPU, with the number of time and/or lens samples a user-specified parameter. Each pass computes a single instant in time and lens position for all pixels. This contrasts with the approach of using a different time and lens position for each pixel or sample [CPC84, Coo86].

The passes can be accumulated either by rendering into a single buffer multiple times, or by rendering into multiple buffers simultaneously (Listing 2). The former reduces memory usage at the expense of repeated traversals of the geometry, while the latter increases memory usage for the buffers, but minimizes geometry traversals and any per-primitive computations and state changes. With either method, the same geometry will be rendered the same number of times by OpenGL. Accumulation is always done on the GPU, to avoid readback.

If the time-sample buffers exceed the available video memory, performance drops radically, as buffers are swapped back and forth to CPU memory. Therefore we prefer to use multiple passes.

Rather than transferring grid data to the video memory for

---

**Listing 2** Pseudo-code for the accumulation algorithms.

---

```

Multibuffer accumulation:
    for every primitive P:
        for every time sample  $t$  from 0 to T-1:
            render P( $t$ ) into buffer[ $t$ ]
    for every time sample  $t$  from 0 to T-1:
        accumulate buffer[ $t$ ] into final image

Multipass accumulation:
    for every time sample  $t$  from 0 to T-1:
        for every primitive P:
            render P( $t$ ) into buffer
        accumulate buffer into final image

```

---

each pass, we cache them in *vertex buffer objects* (VBOs) [SGI03]. This can double motion blur performance. We have not seen swapping issues with VBOs, perhaps because we use them only for grids above a certain size, and one bucket will generally have a limited number of these visible.

Occlusion-culled grids provide another opportunity for performance improvement over CPU culling. Motion-blurred bounding boxes are particularly inefficient for occlusion culling, whereas grid culling at each time sample is easily implemented as part of our motion blur rendering.

##### 5.1. Vertex Motion

Each pass renders the geometry for a specific time in the shutter interval, and in the case of depth of field, a specific position on the lens. We use a GPU vertex program that performs the motion blur interpolation, the model-to-view transformation, and the lens position offset (based on depth; see formula in [PC81]).

Sampling shutter times and lens positions simultaneously for the whole bucket, as opposed to having the time and lens correspondence differ for every pixel, leads to correlated artifacts (strobing, etc.) at low sampling rates. However, with a sufficiently high number of passes (relatively inexpensive with graphics hardware), artifacts become negligible.

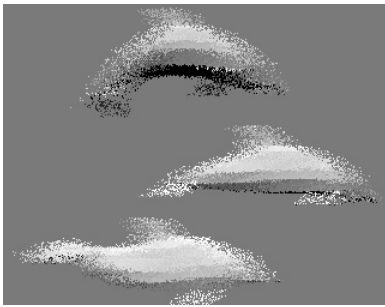
We automatically clamp the number of passes for a bucket based on maximum vertex motion, for motion blur, and maximum area of the circle of confusion, for depth of field. This gives a big speed-up for buckets with little or no blur.

##### 5.2. Per-Pixel Time Sampling

GPUs can also simulate the traditional stochastic sampling technique of associating a specific time value with each sub-pixel sample [CPC84, Coo86]. For each triangle, we rasterize a quadrilateral covering the screen space bounds of the triangle's entire motion, applying a fragment program that (a) interpolates the triangle's vertex coordinates to the time values associated with that pixel, (b) tests the sample point to see if it intersects the interpolated triangle, (c) kills the fragment if it fails the intersection test, and otherwise (d) samples the interpolated triangle to compute the fragment's depth value and color. This results in an image as shown in Figure 6.

Though straightforward, this approach is inefficient for a variety of reasons. The GPU contains dedicated hardware for rasterization, clipping, and culling that is not used by this technique. The moving triangle's bounding box must be computed, either on the CPU or in the GPU vertex program. The resulting rectangular bounds contain many points outside the triangle which must be run through the complex fragment program. The lengthy fragment program that interpolates and tests intersections becomes a bottleneck for GPU throughput. Current GPUs do not efficiently implement *early fragment kill*, which would prevent the computation of the interpolated depth for pixels that fail the hit test. Finally, the fragment program must





**Figure 6:** Using a fragment program to sample a moving geometry at a different time value per pixel. For clarity, we have used just one time sample per pixel.

compute and output  $z$ ; this is called *depth replacement* and diminishes GPU performance. We have found the accumulation techniques we describe to have significantly better performance.

## 6. Filtering and Downsampling

Once a bucket has been completely rendered at the supersampled resolution  $(w_q, h_q)$ , it is downsampled to the final pixel resolution  $(w_p, h_p)$  by a two-pass separable filtering algorithm that runs on the GPU and that supports user-selected high quality filters with large support regions. The downsampled image is then read back to the CPU for output. We avoid any readback of the higher resolution images, since readback is slow.

The first pass convolves the  $x$  filter kernel with the rows of the supersampled image, resulting in an intermediate image of size  $(w_p, h_q)$ . The second pass convolves the  $y$  filter kernel with the columns of the intermediate image.

The fragment program for each pass is generated on the fly by unrolling the user-specified filter into straight-line code that includes all filter weights and pixel offsets as numerical constants. This avoids per-pixel evaluation, or even per-pixel texture lookup, of the filter kernel, and it avoids loop overhead. Details are available in [WE05].

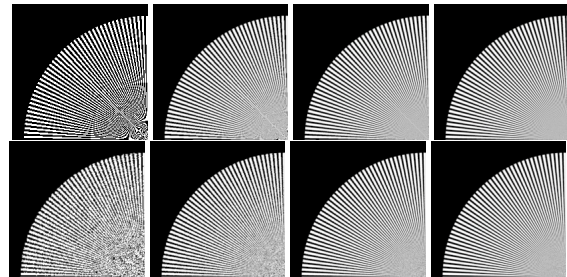
## 7. Results

We have implemented these algorithms in the context of a production-quality renderer (Gelato 1.1). It has been tested on a wide range of scenes, both real-world and contrived, by ourselves and by others. Compared to software-based stochastic sampling renderers, we have found that our algorithm has comparable image quality and superior performance.

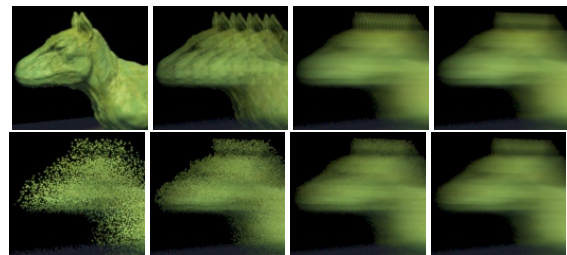
### 7.1. Image Quality

Figure 7 shows a portion of a radial test pattern rendered with our algorithm (top) versus stochastic sampling (bottom) for a variety of sampling rates (from left to right, 1, 4, 16, and 32 samples per pixel). At low sampling densities, the regular sampling of the hardware rasterization shows egregious aliasing, but the superiority of stochastic sampling becomes negligible surprisingly quickly. In real-world examples, noticeable artifacts are even less visible than in pathological examples such as this test pattern.

Figure 8 compares the motion blur of our algorithm's regular sampling (top) with that of stochastic sampling (bottom), for a variety of temporal sampling rates. Below a certain threshold (dependent on the amount of motion in the scene), regular sampling suffers from significant strobing artifacts, while stochastic sampling degrades more gracefully. However,



**Figure 7:** Antialiasing quality with 1, 4, 16, and 32 samples per pixel. Top: regular sampling; bottom: stochastic sampling.



**Figure 8:** Motion blur quality with (from left to right) 1, 4, 16, and 32 temporal samples. Top: regular sampling; bottom: stochastic sampling. Model courtesy of Headus, Inc.

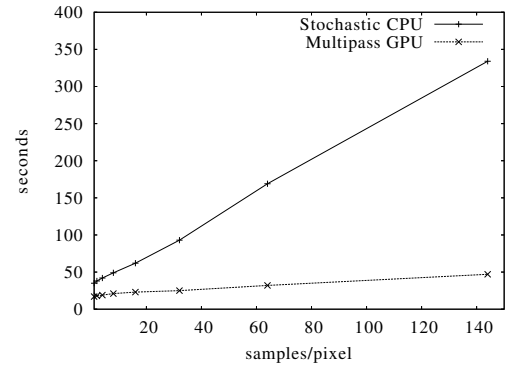
above that threshold, regular sampling gives a smooth appearance, without the grain of stochastic sampling. Regular sampling (both spatial and temporal) will always have visible artifacts or strobing when the sampling rates are not adequate for the scene geometry or motion. Somewhat more care may be necessary for users to choose adequate sampling rates, compared to stochastic sampling. But modern hardware can rasterize at high resolution and with many passes very rapidly, making the necessary sampling rates quite practical, even in a production setting.

### 7.2. Performance

All of the trials discussed below were timed on a 2.1 GHz Athlon 3000 running Linux, with an NVIDIA Quadro FX 3000 (NV35). We report the sum of user and system time, since time spent in the graphics driver can sometimes register as system time. Rendering times are compared against a highly optimized commercial renderer that uses CPU-based stochastic point sampling (PhotoRealistic RenderMan 12.0). The test scene used is either 1 or 8 copies (not instances) of a NURBS character [Headus] of about 17,000 control points, rendered with mapped displacement (not just bump mapping), procedural color, simple plastic shading, and four light sources, one of which is shadow-mapped.

To isolate the time for hidden surface removal, we would like to subtract the cost of the other rendering phases such as geometry management (reading, splitting, dicing) and shading (which includes reading texture and shadow maps from disk). We expect those costs to be fixed with respect to spatial and temporal sampling rates. However, we cannot separate these phases precisely, since hiding informs shading. A low estimate of these fixed costs is the time difference between a full render and one with trivial surface shaders, which we call the *shading delta*.

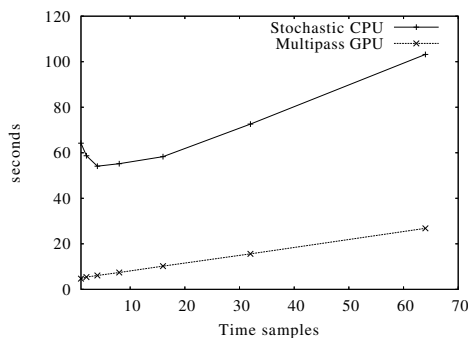
Figure 9 shows rendering time versus spatial sampling for both renderers. Here we have subtracted a measured shading delta of 32 sec. for our GPU-accelerated renderer and 44 sec.



**Figure 9:** Performance comparison of our GPU-assisted multipass algorithm versus a CPU-only stochastic point sampling renderer, rendering the displayed frame at  $1800 \times 1080$  pixels with 1, 2, 4, 8, 16, 32, 64, and 144 samples per pixel. Total render time minus shading delta.

for the CPU renderer. For both renderers, hiding time appears linear in the number of samples, but the marginal cost per sample is 10 times lower for the GPU. At high sampling rates, the GPU hider is much faster.

Figure 10 shows rendering time versus temporal sampling rate. The stochastic CPU renderer renders 16 time samples in about the same time as 1, for a noisy but usable image. At such low sample rates, the GPU renderer is very fast, but is likely to strobe. For an image that requires 32 regular time samples but only 16 stochastic samples, the GPU renderer is still 3 times faster. Each added sample costs about 0.35 sec. on the GPU or 0.94 sec. on the CPU. The stochastic hider has the constraint that adding time samples requires adding spatial samples, because it uses the same samples for both. This is a disadvantage, particularly since, visually, blurrier scenes require fewer spatial samples, not more.



**Figure 10:** Motion blur performance, rendering the full creature shown cropped in Figure 8 at  $1024 \times 786$  pixels. The CPU renderer uses the same samples for time and space; the GPU renderer was run with 16 pixels per sample.

Figure 11 shows a *geometry-heavy* scene, an army of 800 NURBS characters. Rendering it with trivial shaders is over 5 times faster with our GPU renderer than with the CPU renderer. ‡ While it is hard to estimate how much of that time

is hidden-surface removal versus geometry management, it demonstrates that even though large amounts of grid data are being transferred to the GPU, our algorithm still behaves very well.



**Figure 11:** Army of 800 displaced NURBS creatures,  $1800 \times 1100$  pixels, 36 samples per pixel.

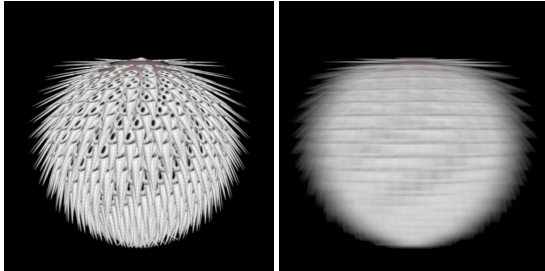
A scene using global illumination or other slow shading methods will be *shading heavy*. Here the hider's speed is less important than how aggressively it culls occluded points before they are shaded. This is more difficult to compare between renderers with differing shading and geometry systems. But as mentioned earlier, the GPU's ability to do relatively fast grid occlusion allows us to cull an extra 10-20% of points, versus a typical Reyes CPU algorithm.

Our GPU renderer slows drastically for scenes with *small grids*, such as a million pixel-sized triangles that each form a one-micropolygon grid. Small batch sizes are the bane of GPU performance. A future project is to reduce this problem by combining nearby grids. We would not expect a CPU hider to have this issue.

The slowest cases for the GPU-based hider are those that require many passes. Modest motion blur performs well, as we've seen. But a bucket containing a grid that moves 100 pixels will require about 100 motion blur passes to avoid strobing. So *extreme motion* gets expensive. Extreme motion increases noise in stochastic hidings, but that is usually visually acceptable. Similar issues apply to depth of field. Scenes with *layered transparency* also require the GPU to rasterize grids multiple times, for depth peeling, whereas the CPU can

‡ Stochastic CPU: 721 sec. full shaders, 554 sec. trivial shaders. Multipass GPU: 163 sec. full shaders, 96 sec. trivial shaders.

maintain arbitrary lists of fragments per pixel. Combining motion blur with transparency multiplies the required number of passes, and eventually the CPU wins. Figure 12 shows a sphere with 550 transparency-mapped “feathers”. Even without motion blur, this example runs 50% slower in our renderer than in the CPU renderer. But in motion with 64 time samples, the stochastic renderer is nearly 9 times faster (stochastic CPU: 7.5 sec. static, 59 sec. moving; Multipass GPU: 12.0 sec. static, 525 sec. moving).



**Figure 12:** A case where our algorithm performs poorly: many stacked motion-blurred transparent objects. The “feathers” are rectangular strips that use a texture to modulate transparency.

## 8. Conclusions

We have described an algorithm for hidden-surface removal that leverages commodity graphics hardware to achieve superior quality and features compared to traditional hardware rendering, while outperforming traditional CPU-software-based high-quality rasterization for typical scenes.

This paper makes the following contributions:

- An algorithm that systematically incorporates high-end features into a hardware-oriented rendering framework. These include supersampling at very high rates, user-selected filters with arbitrarily wide support, motion blur and depth of field, order-independent transparency, multi-channel opacity, and multiple output channels.
- Two optimizations of the depth peeling technique, opacity thresholding and z-batches, that allow it to perform in practice as  $O(N)$  rather than  $O(N^2)$ .
- An exploration of the performance of the algorithm in various cases, and of a variety of engineering trade-offs in using GPUs for high-quality hidden surface removal.

It is often assumed that regular spatial sampling will give inferior results to stochastic sampling. When rates of 4 or 16 samples per pixel were considered high, that may have been true. But our experience has been that at the dozens to hundreds of samples per pixel that are easily affordable when leveraging graphics hardware, the deficiencies of regular sampling are visible only with contrived pathological examples, and not in “real” scenes. If desired, artifacts from regular sampling could be further reduced by using interleaved sampling [KH01] or multisampling (when supported by graphics hardware).

Using advanced GPU features such as floating point precision and detailed occlusion queries can cause current GPU drivers and hardware to run at reduced speed, perhaps 4-16x slower than the optimized paths. Furthermore, despite our algorithm’s minimal use of readback, we find the GPU is idle much of the time. With future work to hide more GPU latency, new GPU designs with fewer penalties for high-precision computation, and the integration of GPU-assisted shading, we hope to recapture this lost performance.

In conclusion, high-quality rendering systems can now be

built on a substrate of commodity graphics hardware, for off-line rendering as well as for real-time or interactive applications. The considerable power of the GPU can be leveraged without compromising either image quality or advanced features. We expect that similar hybrid hardware/software solutions will become more common as GPUs continue to improve in speed and capability.

## References

- [AG99] APODACA A. A., GRITZ L.: *Advanced RenderMan: Creating CGI for Motion Pictures*. Morgan-Kaufmann, 1999.
- [ATI03] ATI: OpenGL Extension Specifications (<http://www.ati.com/developer/>), 2003.
- [CCC87] COOK R. L., CARPENTER L., CATMULL E.: The Reyes image rendering architecture. In *Computer Graphics (SIGGRAPH '87 Proceedings)* (July 1987), vol. 21, pp. 95–102.
- [Coo86] COOK R. L.: Stochastic sampling in computer graphics. *ACM Transactions on Graphics* 5, 1 (Jan. 1986), 51–72.
- [CPC84] COOK R. L., PORTER T., CARPENTER L.: Distributed ray tracing. In *Computer Graphics (SIGGRAPH '84 Proceedings)* (July 1984), vol. 18, pp. 137–45.
- [Eve01] EVERITT C.: *Interactive order-independent transparency*. Tech. rep., NVIDIA Corp. (<http://developer.nvidia.com/>), 2001.
- [GK93] GREENE N., KASS M.: Hierarchical Z-buffer visibility. In *Computer Graphics Proceedings, Annual Conference Series, 1993* (1993), pp. 231–240.
- [HA90] HAEBERLI P. E., AKELEY K.: The accumulation buffer: Hardware support for high-quality rendering. In *Computer Graphics (SIGGRAPH '90 Proceedings)* (Aug. 1990), vol. 24, pp. 309–318.
- [Hea] HEADUS: Killeroo 3D Scan <http://www.headus.com/au/samples/killeroo>.
- [KGP\*94] KELLEY M., GOULD K., PEASE B., WINNER S., YEN A.: Hardware accelerated rendering of csg and transparency. In *Proceedings of SIGGRAPH 94* (July 1994), Computer Graphics Proceedings, Annual Conference Series, pp. 177–184.
- [KH01] KELLER A., HEIDRICH W.: Interleaved sampling. In *Rendering Techniques 2001: 12th Eurographics Workshop on Rendering* (June 2001), pp. 269–276.
- [Kir02] KIRK D.: When will raytracing replace rasterization? SIGGRAPH Panel, 2002.
- [Lev90] LEVOY M.: Efficient ray tracing of volume data. *ACM Transactions on Graphics* 9, 3 (July 1990), 245–261.
- [Mam89] MAMMEN A.: Transparency and antialiasing algorithms implemented with the virtual pixel maps technique. *IEEE Computer Graphics & Applications* 9, 4 (July 1989), 43–55.
- [MH02] MÖLLER T., HAINES E.: *Real-Time Rendering*, second ed. A K Peters, 2002.
- [NVI03] NVIDIA: NVIDIA OpenGL Extension Specifications (<http://developer.nvidia.com/>), Jan 2003.
- [PC81] POTMESIL M., CHAKRAVARTY I.: A lens and aperture camera model for synthetic image generation. In *Computer Graphics (SIGGRAPH '81 Proceedings)* (Aug. 1981), vol. 15, pp. 297–305.
- [PD84] PORTER T., DUFF T.: Compositing digital images. In *Computer Graphics (SIGGRAPH '84 Proceedings)* (July 1984), vol. 18, pp. 253–259.
- [SA03] SEGAL M., AKELEY K.: *The OpenGL Graphics System: A Specification (version 1.5)* <http://www.opengl.org/>, 2003.
- [SGI03] SGI: OpenGL Extension Registry <http://oss.sgi.com/projects/ogl-sample/registry/>, 2003.
- [TK96] TORBORG J., KAJIYA J.: Talisman: Commodity real-time 3d graphics for the pc. In *Proceedings of SIGGRAPH 96* (Aug. 1996), Computer Graphics Proceedings, Annual Conference Series, pp. 353–364.
- [WE05] WEXLER D., ENDERTON E.: *GPU Gems II*. 2005, ch. 21, High-Quality Antialiased Rasterization, pp. 331–344.
- [Wit01] WITTENBRINK C. M.: R-buffer: A pointerless a-buffer hardware architecture. In *2001 SIGGRAPH / Eurographics Workshop on Graphics Hardware* (Aug. 2001), pp. 73–80.