Compressed Coverage Masks for Path Rendering on Mobile GPUs

Pavel Krajcevski, Dinesh Manocha, Fellow, IEEE,

Abstract—We present an algorithm to accelerate resolution independent curve rendering on mobile GPUs. Recent trends in graphics hardware have created a plethora of compressed texture formats specific to GPU manufacturers. However, certain implementations of platform independent path rendering require generating grayscale textures on the CPU containing the extent that each pixel is covered by the curve. In this paper, we demonstrate that generating a compressed grayscale texture prior to uploading it to the GPU creates faster rendering times in addition to the memory savings. We implement a real-time compression technique for coverage masks and compare our results against the GPU-based implementation of the highly optimized Skia rendering library. We also analyze the worst case properties of our compression algorithms. We observe up to a 2X speed improvement over the existing GPU-based methods in addition to up to a 9:1 improvement in GPU memory gains. We demonstrate the performance on multiple mobile platforms.

Index Terms—texture compression, coverage masks, 2D path rendering.

1 INTRODUCTION

O NE of the main challenges in computer graphics is play objects at a finite resolution. Improper discretization may lead to noticeable aliasing artifacts due to insufficient sampling. In order to alleviate these artifacts, different techniques have emerged for computing proper discretizations [1][2]. When rasterizing geometric objects, the main difficulty is determining what percentage of a pixel is covered by the screen-space projection of the object. This information, once calculated, can be stored in an image known as a *coverage mask*. Coverage masks are usually stored as eight-bit grayscale images and can be used in a variety of different ways in order to speed up the rendering of geometric primitives, including caching [3] and GPU based rendering of 2D curves [4].

Pixel coverage remains an instrumental part of proper rasterization. There are many applications where coverage masks are useful, from culling [5] to visibility determination for more efficient lighting [6]. In this paper, we mainly focus on coverage masks used in rendering non-convex piece-wise two-dimensional cubic and quadratic curves, or *paths*, with anti-aliasing (Figure 1). These curves are used in a majority of vector graphics data, most importantly as the basis for resolution-independent text rendering using different fonts and sizes. These coverage masks, generated at run-time from network data such as web pages, are used billions of times on a daily basis [7]. To further motivate the problem, we have traced the rendering procedures of over 750,000 web pages from the Chrome internet browser. Of these web pages, we observed that 51% draw arbitrary paths of which 19% are anti-aliased requiring dynamically generated textures. Of the paths that require coverage information,

E-mail: pavel@cs.unc.edu and dm@cs.unc.edu

Manuscript received -; revised -



Fig. 1: (*Top left*) The piece-wise anti-aliased cubic curve used as input. (*Bottom Left*) The final rendered curve. (*Top right*) The uncompressed coverage mask passed to the GPU to determine the amount each pixel is covered by the curve. (*Bottom right*) The compressed coverage mask using our method. On the far right is a zoomed in comparison of the compressed and uncompressed masks. Although only a few pixels differ, using our method, these masks are compressed in real time and save time and memory during the rasterization of these curves.

most of the web page rendering time is spent drawing the coverage mask of the path on the CPU prior to uploading it to the GPU.

In this paper, we show that coverage masks generated at run-time by the CPU can be compressed efficiently for GPU-based rendering with little loss in rendering fidelity. We present a way to augment the scan conversion process

P. Krajcevski and D. Manocha are with the Department of Computer Science, University of North Carolina at Chapel Hill, Chapel Hill, NC, 27599.

of non-convex path rendering to directly output compressed textures for use on GPUs with corresponding texturing hardware support. We demonstrate encoding into a variety of different compression formats in order to show applicability to a widespread range of commodity graphics hardware. In particular, we show that even with general 32-bit CPUs, efficient coverage mask compression can be performed to target the widely used DXTn, ETC, and ASTC texture compression formats [8][9][10]. Finally, we demonstrate a speedup of up to 2X in rendering performance using compressed coverage masks on current mobile platforms (e.g. tablets and smart phones). This savings in rendering speed is in addition to the GPU memory gains of 2X up to 9X depending on the texture compression format. Our method is integrated into the Skia¹ two-dimensional rendering library [4]. This library is the rendering backbone in the popular Google Chrome and Mozilla Firefox web browsers that currently boast billions of users [7].

Additionally, we perform an in-depth analysis of the compression quality of different texture compression formats. We demonstrate worst-case scenarios with respect to texture fidelity and discover that our method meets the requirements to compress coverage masks, yet performs quite poorly for general grayscale data. However, due to the predictable appearance of coverage masks, we can exploit many of their properties to create perceptibly identical renderings on general purpose hardware. The Skia library contains a suite of performance and correctness tests covering both test data and web-page data. Overall, our approach aligns well with the current hardware and software trends. The mobile GPU market is growing at a considerable rate with more than a billion sales per year [11]. To address this trend and develop higher performance on mobile GPUs, hardware vendors are developing more aggressive compression formats that are designed specifically for such GPUs [10]. In particular, energy savings during rendering are becoming more important. Using a few extra CPU operations in order to decrease the texture bandwidth by 2-3X likely produces significant energy savings for texture-heavy mobile applications. Texture memory accesses are almost three orders of magnitude more expensive than standard ALU operations [11]. Our method for compressing coverage masks leverages these trends and becomes increasingly useful with the current architectural trends of modern GPUs.

The rest of the paper is organized as follows. Section 2 gives an overview of recent work in coverage masks and compression formats. Section 3 presents our scan conversion algorithm used during rasterization, and the various compression formats used to store grayscale coverage information. In Section 4 we analyze the compression methods analytically and show scenarios with worst-case compression quality. We highlight the performance of our algorithm on various mobile devices in Section 5. Finally, we present conclusions, limitations, and future work in Section 6. A preliminary version of this paper appeard in [12].

2 BACKGROUND

In this section, we give a brief overview of prior work on coverage masks, GPU-based vector graphics, and texture



Fig. 2: A piece-wise quadratic curve is filled with green using the Loop-Blinn method [13]. The pixels (pink) whose centers are not covered by the triangles circumscribing the curve will not be drawn if the GPU is not using a hardware anti-aliasing method. For power constrained GPUs, such as those on mobile devices, multi-sample anti-aliasing is prohibitively expensive due to the large number of fragment shader invocations. When the curve is non-convex, hardware rasterization tends to generate more inaccurate pixel coverage than software rendering.

compression.

2.1 Coverage Masks

One of the major problems in computer graphics has been to determine the amount that a geometric shape, commonly a triangle, covers a given pixel during rasterization [1][3]. This problem, also known as *pixel coverage*, is used to reduce aliasing artifacts caused by the discrete nature of our display devices and memory layouts. More recently, coverage masks have been used for more than simply anti-aliased rasterization. Zhang et. al. [5] use occlusion maps, a variation of coverage masks, to quickly cull non-visible geometric primitives during the rendering of large scenes. Kautz et. al. [6] use coverage masks to cache hemispherical visibility information in order to perform efficient self-shadowing of objects. Coverage information has also been used to accelerate shading operations in the GPU pipeline, although these methods are more suited to hardware implementations than software

Coverage masks are used extensively to render 2D images from geometric primitives. In particular, coverage information is necessary when rasterizing anti-aliased polygons independent of the color and shading information. In order to render these polygons, first the pixel coverage mask is generated, and then the color of the polygon is modulated by the intensity of the pixel in the coverage mask. This technique is used in the 2D rendering library Skia [4] for GPU rasterization of non-convex anti-aliased paths.

2.2 GPU-based Vector Graphics

Resolution-independent rendering is important for many objects in graphics such as the arbitrary cubic and quadratic curves used to represent shapes in most modern fonts. Until recently, these curves have been rendered using software rasterization algorithms. Given the recent advances in GPU

^{1.} https://sites.google.com/site/skiadocs/home

development, there has been considerable groundbreaking work to use GPUs to perform resolution-independent rasterization [13][14][15]. As pioneers in this work, Loop and Blinn [13] devised a method to rasterize Bézier curves by assigning values to the texture coordinates of triangles derived from the control points of the curve. These values were used to calculate the distance from the curve in the given triangle, which was used for proper anti-aliasing. Kokojima [16] improved the efficiency of this method by exploiting the stencil buffer. Qin [15] presented a method to exploit the texture storage of a graphics processor to store curve information using approximate circular arcs. Finally, Kilgard and Bolz [14] described an approach that transmits control points directly to the GPU to render the curve. Although this method renders vector graphics very quickly, it requires proprietary hardware features such as specific library extensions. Further approaches using signed distance fields have been used by Green [17] for artist generated vector graphics.

2.3 Anti-Aliasing Non-Convex Curves

Despite recent advances in using GPUs to accelerate vector graphics rasterization, certain classes of vector graphics still remain slow on mobile hardware [18]. Of the techniques mentioned in Section 2.2, the Loop-Blinn method is among the fastest techniques for rendering resolution-independent vector graphics from arbitrary path data. The GPU-based method introduced by Kilgard and Bolz [14] builds upon the Loop-Blinn method by implementing a conservative approach to determining coverage information in hardware. Most notably, as shown in Figure 2, for paths that generate smooth curves but are comprised of multiple control points, the triangles that conjoin quadratic and cubic pieces of a curve may not cover all necessary pixels. When these triangles are rasterized by the GPU, the centers of some pixels covered by the path may not be covered by the triangles. For GPUs that do not support hardware-based anti-aliasing, or where such anti-aliasing is too expensive due to power constraints, pixels that should have partial coverage from the path will not be drawn. This can cause aliasing artifacts when rendering curves whose details are on the order of a single pixel.

To support many different use-cases, the 2D rendering library Skia chooses different rendering paths dependent on the path being rendered. For non-convex paths without antialiasing, Skia approximates a path using line segments and then uses their endpoints as input to a triangle fan drawing both front and back facing triangles. Using the stencil buffer, pixels can be turned on or off based on whether they are inside or outside the path. However, line segments create significant aliasing artifacts during rendering, and this technique cannot be used for anti-aliased paths.

To perform anti-aliasing, in certain cases Skia uses the Blinn-Phong method followed by extruding the triangles along the normal to the path by the amount required to cover all of the pixels covered by the path. However, for general non-convex paths, this results in artifacts in areas where the extruded polygons of two different curves overlap leading to double-blending and incorrect pixel coverage. As a result, the GPU-based renderer in Skia draws the coverage information in software prior to uploading the resulting



Fig. 3: The different stages in GPU-based rendering of filled 2D regions using coverage masks. The only part that takes place on the GPU is the compositing. Our contribution in this modified pipelineis the stage outlined in red, where compressed textures are generated directly from the runlength encoded coverage information. In doing so, we avoid both writing a full resolution texture into CPU memory and uploading a full resolution texture to GPU memory, providing savings on both ends.

grayscale texture to the GPU for shading. This rendering algorithm used to support the use of GPUs can become a significant bottleneck during the rendering of anti-aliased concave paths [4]. In this paper, we show that the grayscale coverage information can be efficiently compressed to a texture format (Section 2.4) thereby significantly increasing the speed at which it is uploaded to the GPU.

2.4 Texture Compression Formats

Over the past few decades, there has been significant research into texture representations in GPU memory. The main requirements for texture representation formats were outlined by Beers et al. [19] as random access and hardwarebased decompression. Real-time decoding is supported in modern GPUs, though the performance of the encoding step can be slow and are generally not done in realtime [20]. Over the years, many new compression formats have emerged offering quality versus performance tradeoffs [8][9][10].

One of the earliest texture compression formats introduced in commodity graphics hardware was the DXTn family of compression formats [8]. Variations of this format have been implemented in hardware to support grayscale textures and textures with alpha. Subsequently, Ström and Akenine-Moller introduced ETC1, a texture compression format that uses scale and offset factors from look-up tables to reconstruct pixel values [21]. A few years later, Ström and Petterson introduced ETC2, which improved upon ETC1 by allowing invalid bit combinations to encode a wider range of pixel values [9]. Single channel variations have also been introduced, but their adoption has not reached commodity graphics hardware [22]. Nystad et al. [10] recently unveiled ASTC, which allows encoders to choose between a variety of compression methods and a variable bitrate from eight bits per pixel down to 0.89 bits per pixel. Although this flexibility in the compression format allows a large quality versus compression size trade-off, developing real-time encoders for ASTC can be challenging.

3 COMPRESSED SCAN CONVERSION

In this section we describe our technique for encoding the coverage information into a GPU-based compressed texture format. Given a piece-wise two-dimensional curve, or *path*, we augment the scan conversion algorithm on the CPU for generating coverage information. Our formulation is based on the assumption that the time spent writing the encoded coverage information into a GPU-specific format can be recovered during the time it takes to upload the texture to the GPU. Even if the time saved by uploading a compressed representation is lost during the encoding step, we still gain memory savings from using compressed textures.

The input to our algorithm is a list of 2D curves defined using Bézier control points. From this list, our goal is to generate an accurate two-dimensional grid of pixels that best approximate the curve along with a specified *paint*. The paint determines the color and opacity of the pixels that are covered by the curve along with any other special operations such as anti-aliasing and gradient dithering. For pixels that are partially covered, they will be painted proportional to the amount that they are covered by the path. In a GPU-based rasterization pipeline, the coverage information is first generated and then used as a texture along with the paint to write to the framebuffer.

There are two operations commonly used for rasterizing these paths. First, the path may be *filled* such that a single color is painted within the bounds defined by the path. In this case, the coverage information in conjunction with the paint opacity is used to determine how much of that color should be blended with the background color. If the path is being rendered using the GPU, the coverage information must be uploaded as a texture prior to determining the final color and blending. The other operation, known as stroking, draws an outline of a given thickness along the path. In this case, the Skia library computes a new path along the outline of the stroke. Rendering this new path filled with the stroke color is identical to rendering the original stroked path. We restrict our formulation to non-convex paths. Convex paths can be efficiently drawn on GPUs by using a triangle fan in conjunction with the stencil buffer in a modified Loop-Blinn method described in Section 2.2 [4].

The texture uploaded to the GPU is the image that stores the pixel coverage information. We proceed by first describing a variety of compression methods that we use to encode grayscale information on commodity graphics hardware. We then describe how we augment the scan conversion process to rows of compressed texture data.

3.1 Compression Formats

Due to the large schism of hardware support for various texture compression formats, our goal is to develop an approach that is portable between different GPUs. Decoding algorithms tend to be relatively simple because of the necessity of hardware-based implementations of GPU-encoded textures. Our encoding algorithm exploits this simplicity inherent in all compression formats. As described in Section 3.2, neighborhoods of pixels in coverage masks usually contain either fully transparent or fully opaque pixels. This allows us to precompute many of the parameters for our compression formats prior to the actual encoding. However, the reconstruction of the coverage information from these formats is necessarily lossy, due to the nature of the random access constraints. The following is a detailed overview

```
uint32_t BytesToDXTnIndices(uint32_t x) {
    // Collect and invert high three bits
    x = 0x07070707 - ((x >> 5) & 0x07070707);
    // Set mask if any bits are set
    const uint32_t mask = x | (x >> 1) | (x >> 2);
    // Mapping: 7 6 5 4 3 2 1 0 -> 8 7 6 5 4 3 2 0
    x += mask & 0x01010101;
    // Handle overflow:
    // 8 6 5 4 3 2 1 0 -> 9 7 6 5 4 3 2 0
    x |= (x >> 3) & 0x01010101;
    // Result: 9 7 6 5 4 3 2 0 -> 1 7 6 5 4 3 2 0
    return x & 0x07070707;
}
```

Fig. 4: C code for converting an integer storing four 8bit values into four three-bit indices corresponding to the proper layout of a DXTn block. Using branchless code without multiplies or divides yields extremely fast and pipelined code on modern CPU architectures.

```
uint32_t BytesToETC2Indices(uint32_t x) {
  // Three high bits: 0 1 2 3 4 5 6 7
  x = 0x07070707 - ((x >> 5) \& 0x07070707);
// Negate: 0 -1 -2 -3 -4 -5 -6 -7
  x = ((0x80808080 - x) ^ 0x7F7F7F7F);
  // Add three: 3 2 1 0 -1 -2 -3 -4
  const uint32_t s = (x \& 0x7F7F7F7F) + 0x03030303;
  x = ((x ^ 0x03030303) & 0x80808080) ^ s;
  // Absolute value...
  const uint32_t a = x & 0x80808080;
const uint32_t b = a >> 7;
  // M is three if the byte was negative
  const uint32_t m = (a >> 6) | b;
  // .. continue absolute value:
  // 3 2 1 0 1 2 3 4
  x = (x \hat{} ((a - b) | a)) + b;
  // Add three to the negatives:
  // 3 2 1 0 4 5 6 7
  return x + m;
}
```

Fig. 5: C code for converting an integer storing four 8bit values into four three-bit indices corresponding to the proper layout of an ETC2 block. Similar to Figure 4, we perform the conversion using only bitwise operations and without expensive multiplies or divides.

of the algorithm applied to the DXTn, ETC2, and ASTC families of compression formats.

3.1.1 DXTn

In the DXT family of texture compression formats, introduced by Iourcha et. al. [8], 4×4 pixel blocks are encoded by storing two pixel values per block and a twobit index per pixel. The two separate pixel values stored in the block generate a palette of colors from which the perpixel index selects the final color. The palette is based on intermediate values chosen by linearly interpolating the two stored pixels. For coverage information, we use the DXTn format designed specifically for grayscale known as LATC, or Luminance-Alpha Texture Compression (also known as RGTC, 3DC, and BC4). This format supports two eight-bit grayscale values and sixteen three-bit index values per pixel for a total of 64 bits per block, giving a compression ratio of two-to-one for grayscale images. In order to reach the full range of grayscale values, we store 0 and 255 as endpoints for our coverage mask. Due to the indexing scheme of DXTn, the mapping of coverage values to interpolation indices can not be directly copied from the high three bits of each coverage value. We first quantize each grayscale value to three bits such that their reconstruction into eight bits by bit replication minimizes the error from the original grayscale value. Once these three bits are computed, we must use a mapping from the quantized bits to the proper DXTn indices

$$0, 1, 2, 3, 4, 5, 6, 7 \rightarrow 1, 7, 6, 5, 4, 3, 2, 0.$$

This mapping can be performed without branches on commodity hardware using eight bits per index. If we treat each block row as four 8-bit grayscale values, we can store an entire block row in a single 32-bit register. Furthermore, 32-bit integer operations can be used to perform byte-wise SIMD computations without requiring special SIMD hardware, as shown in Figure 4.

3.1.2 ETC2

One variant of the ETC2 compression format is a table-based compression algorithm that takes 4×4 blocks of grayscale pixels, and reconstructs 11-bit grayscale values from 64-bit encoded data in order to provide higher precision than traditional 8-bit textures. However, the 64-bit representation maintains a two-to-one compression ratio similar to DXTn. The procedure by which the coverage value for pixel c_i is reconstructed is

$$c_i = b \times 8 + 4 + (\mathbf{T}_v)_{t_i} \times 8,$$

where the encoded data stores an 8-bit base codeword b, a 4-bit multiplier m, a 4-bit modulation index v, and sixteen 3-bit indices t_i . **T** is a table containing sets of modulation values constant across all the encodings. This table has sixteen entries, indexed by v. Each t_i selects a final modulation value from the set \mathbf{T}_v . The result c_i is then clamped to the range [0, 2047].

To compress the grayscale coverage information, we first fix values for v, b, and m such that they generate the tightest bounds to the entire range of grayscale values. We compute these values by performing an exhaustive search through all possible combinations of v, b, and m offline. In order to compress the coverage information, we perform a quantization to three bits as described in Section 3.1.1. However, due to the indexing method of ETC2, we must use a different mapping

$$0, 1, 2, 3, 4, 5, 6, 7 \rightarrow 3, 2, 1, 0, 4, 5, 6, 7.$$

This mapping is also has the same implementation advantages as DXTn, as shown in Figure 5, allowing branchless computation to be done in fixed 32-bit registers.

3.1.3 ASTC

Finally, we demonstrate fast compression of our coverage information using the ASTC format introduced by Nystad et. al [10]. This format has a variable block size that must be chosen prior to compression, and we have noticed that even at the highest compression rate, 12×12 , rendering artifacts

were negligible. This is possible due to the high compressibility resulting from the low entropy of the coverage mask described in Section 3.2.

ASTC encoded blocks may choose from many different compression options. One such option is whether or not to partition the block into separate subsets of pixels with different compression parameters. Similar to DXTn and ETC2, ASTC uses per-pixel indices to reconstruct the block of pixels. However, there may be fewer indices than pixels, in which case the indices are stored in a grid and interpolated across the block. Finally, similar to DXTn, ASTC reconstructs pixels by using generated indices to lookup palette entries. However, ASTC allows the block encoding to choose how many bits are allocated towards endpoint representation versus index representation.

In order to maximize the fidelity of the ASTC compressed coverage mask, we outline a list of the choices that we made for each 12×12 block of pixels. The main insight is to maximize the number of pixel index values and their bit depth. We are able to maximize the index size because the endpoints must cover the full range of grayscale values and hence require very few bits. For this reason, we are able to generate a valid ASTC encoding using the following choices:

- 6 × 5 texel index grid to maximize the number of samples in a 12 × 12 pixel block
- Three bits per texel index
- Single plane encoding (redundant due to singlechannel input). This is chosen because we do not use multi-channel pixels
- Only one color endpoint mode: direct luminance
- Single partition encoding with two 8-bit endpoints: 0, 255

Using these constants for all coverage information, there is no special need for the base-three and base-five integer sequences supported by ASTC [10]. Since we know the dimensions of the grid versus the dimensions of the block size, we can precompute the amount that each pixel contributes to each index, and store this in a look-up table. During compression, for each texel grid index we store the top three bits of a weighted average of the pixels that are affected by the index. The final result is 144 grayscale pixels compressed into 128 bits, providing a compression ratio of nine to one. Although compression of ASTC is slower than DXTn and ETC2, the generated compressed textures are significantly faster to load into GPU memory.

3.2 Scan conversion

While the compression format chosen is dependent on the underlying hardware, the scan conversion of path data is computed independently on the CPU. In particular there are two main steps:

- Determine the run-length encoded coverage information for each scanline of pixels
- 2) Convert multiple scanlines at once into the necessary compression format

From a given path, coverage information for each pixel is computed by sampling the path N times per pixel,

						• • • • • • • • • • • •	•	••••	
1	2	2	17						
64	225	192	225						
1	2	2	2	1	35				
0	96	0	127	112	0				

Fig. 6: Sparse run length encoded (RLE) buffers. These buffers are used to store the coverage information for a row of pixels prior to writing them into the coverage mask. For each pixel row, the RLE buffer is allocated to contain as many RLE entries as there are pixels. The scan converter operates on rows of super-sampled pixels, shown here as a 4×4 grid within each pixel, and updates the corresponding RLE buffer. In this figure, the blue entries contain the number of runs of the corresponding pixel value. Grey entries are uninitialized and never written to nor read. Samples which contribute to the coverage of the red curve are drawn in blue and samples that are uncovered are drawn in black.



Fig. 7: Our scan conversion pipeline augmented to output GPU-compressed blocks. For $M \times M$ compressed block sizes, our pipeline operates on M sparse RLE buffers in parallel (Figure 6). Once M columns are processed, they are compressed into the target compressed format. For a given column, we read from the entries in the associated sparse RLE buffers. If any of the row values have changed, we update the corresponding pixel for the current column (outlined in red). Otherwise, we simply copy the previous column. For 8-bit coverage values and 4x4 compressed block sizes, each column fits in a single 32-bit register.

commonly N = 16 with the samples arranged in a regular grid (Figure 6). Each sample is applied a boolean value $b_i \in \{0, 1\}$ such that the final coverage for a given pixel in image I is

$$\mathbf{I}(x,y) = \frac{1}{N} \sum_{i=1}^{N} b_i.$$

For a value corresponding to N = 16, this implies that I can take up to 17 possible values for any $(x, y) \in \mathbf{N} \times \mathbf{N}$.

In a scanline of samples, the edges of the curve can be computed analytically in order to properly set the corresponding b_i . As shown in Figure 6, the per-pixel coverage information, i.e. the number of samples covered by the path, is stored in a sparse run-length encoded (RLE) buffer. This buffer is updated for each new scanline of samples within a row of pixels. The sparsity of the buffer prevents unnecessary allocation when an initial scanline of samples is altered by a subsequent scanline. In this situation, the samples within a pixel may be identical in the first scanline of samples but different in the second.

The pixels containing intermediate values, i.e. those that are neither fully opaque (covered) or transparent (uncovered), are only found along the boundaries of the 2D path. For this reason, a majority of the pixels in a coverage mask take extremal values (0 or 255) and very few, along the edges of the path, tend to have intermediate values. This means that most of the image can be stored as a binary image, producing an entropy close to one [23]. This extremely low entropy property of coverage masks makes them highly compressible.

In order to generate compressed textures, we must adhere to the random access requirements in texture representations. Random access ensures that the renderer has equal access to all pixels regardless of when they are needed. This requirement implies a fixed block size for each compression format: 4×4 for DXTn and ETC, and 12×12 for ASTC. Once a scanline of pixels is computed, it can be stored in a row of an 8-bit grayscale texture. We generate compressed representations of the grayscale textures by consuming M rows of run-length encoded data at a time, where M is the dimension of the (square) block size of the texture compression format. As shown in Figure 7, we read the leftmost column of grayscale values and update the corresponding byte as we walk down our M RLE buffers. At each step, we advance to the column with the earliest ending run length. Once we advance past M columns, we efficiently compute a compressed representation of the $M \times M$ block that we have read from the RLE buffers, as described in Section 3.1. For the most common case, M = 4, the four grayscale values are represented as a 32-bit integer, and we can perform SIMD byte-wise operations using integer shifts and adds. As an optimization, if we advance the current column farther than *M* pixels at once due to the RLE encoding, we can copy the previous block encoding into its neighbor to the right.

4 ERROR ANALYSIS

The methods for compressing coverage masks outlined in Section 3 are designed for speed and with the assumption that coverage masks will be mostly coherent. For any given coverage mask, the rendering time will be dependent on the resolution of the coverage mask. However, the quality is fixed due to the precomputed compression parameters for each format. As a result, it is possible to find the worst-case texture quality for data compressed into each format. In this section we investigate such failure cases and show scenarios that our method might not handle particularly well.

Due to the nature of coverage masks, the only areas of high detail are border regions where pixels may end up being partially covered. As a result, the error bounds reported in this section do not reflect the quality of the final compressed coverage masks. They are used to demonstrate the limitations of our compression scheme against generalpurpose data and highly incoherent coverage masks. In practice, compressed coverage masks do not contain any



Fig. 8: Quantization error when converting the incoming number of samples covered per pixel to the final value stored in the compressed format. We show absolute error for both DXT and ETC formats with respect to the original quantized values. For fully opaque and fully transparent pixels we have no error as designed. For intermediate values, discrepancies in error arise from the way values are quantized in adherence to the two texture formats.

visible artifacts. In contrast to the artifacts that are most commonly noticeable in low-resolution coverage masks, such as aliasing, or "jaggies", the most noticeable artifacts in compressed coverage masks tend to be blurring caused by the interpolation described in Section 4.2.

4.1 DXTn and ETC2 Compression Formats

In both DXT and ETC2, we generate a fixed color palette into which we compress our coverage masks. For both formats, the palette is precomputed based on what the anticipated data in the block will be. Our compression parameters are chosen such that we represent values ranging from fully transparent, or zero, to fully opaque, or $2^8 - 1$.

As described in Section 3.2, our input texture has at most 17 values, ranging from zero to sixteen, which counts the number of samples covered by our path. When uploading uncompressed coverage masks, each of these seventeen values gets quantized to a value from zero to $2^8 - 1$. However, since both DXT and ETC2 use three-bit indices, each compressed block contains only eight possible choices which vary depending on the format. For DXT, the available values are

$$\{0, 36, 73, 109, 146, 182, 219, 255\}$$

while for ETC2, the values are

 $\{0, 51, 78, 105, 149, 176, 203, 255\}.$

Figure 8 shows the amount of absolute error each of the original 17 values incurs when compressing to the respective formats.

4.2 ASTC Compression Format

Unlike DXT and ETC, ASTC blocks interpolate their indices from a low-resolution index grid to determine per-pixel index values. For this reason, determining the proper ASTC representation requires more processing than converting pixels to index values. As in Figure 9, each pixel in the input



Fig. 9: For a 12x12 ASTC block, we maximize the number of samples we store in order to get the finest granularity of control possible over the resulting pixels. Physical limitations of the ASTC format restrict us to a 6x5 index grid stored on disk (red samples). During decompression, these indices are interpolated to each texel (blue samples) to compute the final index used for selecting from the precomputed palette.



Fig. 10: (Top row) Uncompressed failure cases for certain 12x12 blocks. (Bottom row) Our ASTC compression method applied to each block. Due to the interpolation of index coordinates in ASTC blocks, certain blocks will be compressed much more poorly than others. In particular, blocks that have many uncorrelated neighboring pixels, while able to be represented using ASTC, are not particularly well suited for our method. However, such blocks are very rare in coverage mask textures.

block contributes to the final value of four surrounding indices.

In order to maintain real-time performance of coverage mask compression, we must precompute many of the parameters for each block, as described in Section 3.1.3. This optimization has implications on texture compression quality when dealing with high-frequency data. In particular, data that has a large variance between our index locations can become distorted. As we can see in Figure 10, blocks that have high frequency are very difficult to compress using our chosen parameters. In particular, the checkerboard block, which is simply alternating black and white pixels, results in the most artifacts due to high index averaging of all nearby pixel values.

In order to properly select the indices for our ASTC block, we may solve a linear system of the form $\mathbf{A}x = b$, where \mathbf{A} is a 144×30 matrix corresponding to the contribution of each pixel in a 12×12 block to each index in a

Percentage Improvement in Rasterization Speed



Fig. 11: Performance improvements using compressed textures on a variety of different benchmarks. Two of the tests performed were on tablet versions of popular websites. The Google Spreadsheets benchmark data was gathered from the desktop version of the site using many stroked paths. The other two were the vector images in Figure 12.

 6×5 grid. By virtue of **A** being fixed, we can precompute the pseudo-inverse $\mathbf{M} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ in order to find the index values

$$x \approx \mathbf{M}b$$

for any block *b*. Furthermore, we can use this method to determine what the error is for any given block *b*,

$$E(b) = \|\mathbf{A}\mathbf{M}b - b\|_2.$$

We use a least-squares formulation in order to minimize the appearance of noisy pixel values. However, this error function is highly non-convex due to the nature of the quantization of each valid value of b. For this reason, we cannot analytically derive a global maximum or minimum. In pursuit of a numerical solution, we can calculate the gradient for E(b),

$$\nabla E(b) = \left\| \hat{\mathbf{M}} b - b \right\|_2 \left(\hat{\mathbf{M}}^T \hat{\mathbf{M}} b - (\hat{\mathbf{M}} + \hat{\mathbf{M}}^T) b - b \right),$$

where $\mathbf{M} = \mathbf{A}\mathbf{M}$. Using this gradient, we use gradient descent to find the worst-case blocks that can be encoded with our method. Due to the large number of local maxima within the search space, we seed our optimization routine with random blocks. The resulting block can be seen in Figure 9.

The error analysis for ASTC blocks allows us to quickly determine whether or not a given block is suited for compression. If compressing the block will introduce a significant amount of unacceptable compression error, we may abort the compression procedure and try alternatives such as a different format or reverting to uncompressed textures. Additionally, this technique of determining error can aid content authors in creating paths that can be compressed well at various resolutions.

5 RESULTS

To test our results, we have integrated our real-time compression pipeline into the 2D graphics library Skia [4]. This



Chalk	358ms	370ms	371ms	473ms	5ms				
Car	368ms	385ms	385ms	403ms	2ms				
Crown	121ms	127ms	137ms	200ms	15ms				
Dragon	92ms	94.3ms	96ms	140ms	7ms				
Polygon	149ms	152ms	154ms	208ms	5ms				
Compressed									
Image	Min	Median	Mean	Max	σ				
Tiger	81ms	83ms	83ms	93ms	2ms				
Chalk	339ms	349ms	350ms	495ms	5ms				
Car	364ms	387ms	386ms	424ms	3ms				
Crown	106ms	109ms	111ms	168ms	9ms				
Dragon	87ms	92.2ms	101ms	156ms	19ms				
Polygon	133ms	134ms	137ms	194ms	7ms				

Fig. 12: Rendering times of the following images on a first generation Moto X (1.7 GHz Qualcomm Krait, Qualcomm Adreno 320) from 100 runs. From left to right the images are labeled Tiger, Chalk, Car, Crown, Dragon, Polygon.

library is used as the backbone to many cross-platform 2D programs and operating systems including Android, Google Chrome, and Mozilla Firefox. In order to maintain performance and regression tests across all platforms, Skia includes two types of comprehensive tests. For any given change to the implementation, Skia tests the new rendered image against existing baseline images. If any pixels differ by a significant amount, these tests fail and the change is invalid. The second test measures performance against a suite of microbenchmarks and a suite of rendering commands that are invoked during the rendering of common web pages. In order for these tests to pass, their running time must be within a small threshold of the previously passed test. In each of our examples, we have maintained both correctness and performant code with respect to the existing implementations.

First, we must show that our implementation runs fast on modern hardware. In Figure 11, we show different classes of benchmarks that have been run on a variety of different mobile GPUs. In each case, we see a general increase in the rendering speed of certain web pages and common vector graphics benchmarks. As we can see, the desktop GPU does not receive as much of a benefit from the compression routine as the mobile GPUs. We conjecture that mobile GPUs are more sensitive to transmitting large amounts of data from the CPU to the GPU due to power

Mobile Platform	СРИ	GPU	Uncompressed	Compressed	Texture Format	Memory Benefit
Moto X	1.7 GHz Qualcomm Krait	Qualcomm Adreno 320	163ms	137ms	ETC2	2:1
Galaxy Note 3	1.9 GHz ARM Cortex-A15	ARM Mali-T628	171ms	161ms	ETC2	2:1
HTC One M8	2.3 GHz Qualcomm Krait 400	Qualcomm Adreno 330	114ms	102ms	ETC2	2:1
Galaxy Note 10.1	1.9 GHz ARM Cortex-A15	ARM Mali-T628	171ms	136ms	ASTC	9:1
Galaxy S5	1.3 GHz ARM Cortex-A7	ARM Mali-T628	311ms	157ms	ASTC	9:1

TABLE 1: The rendering times for the polygon benchmark (Figure 12) from Skia using both compressed and uncompressed texturing on a variety of CPU/GPU combinations. The polygon benchmark generates a large sequence of thin, concave polygons and stores them as piece-wise 2D paths on the GPU. These polygons are then both stroked and filled to generate a large amount of paths that must be rasterized. From these results, we notice an increase in rendering speed of the heavily optimized Skia library on all mobile devices. Most importantly, the increase in memory efficiency from ETC2 (2:1 ratio) to ASTC (9:1 ratio) provides significant improvements in rendering time. These results were generated from the mean runtime of 100 executions.

restrictions and hence receive more benefits. Mobile GPU performance increases are better demonstrated in Table 1 where various mobile GPUs render the polygon image (Figure 12) from the Skia performance tests. From this table, we observe that both CPU speed (Galaxy Note 10.1 vs Galaxy S5) and compression ratio (Galaxy Note 10.1 vs Galaxy Note 3) play a vital role in rendering performance on mobile devices.

In order to test accuracy, we perform both a visual comparison against the reference images (without compression) and measure the difference using the *Peak Signal to Noise Ratio*, or PSNR:

$$PSNR = 10 \log_{10} \left(\frac{3 \times 255^2 \times w \times h}{\sum_{x,y} \left(\Delta R_{xy}^2 + \Delta G_{xy}^2 + \Delta B_{xy}^2 \right)} \right)$$

In Figure 13, we compare the various use cases of rendered paths and the difference in their rendering. We observe that only pixels along the borders of the paths are affected by the compression scheme. This homogeneity in the coverage masks is the primary reason why they are highly compressible. From the zoomed in comparisons, we notice that there is little to no quality loss in the final images. However, the pixels that differ do so by a non-trivial amount. This difference causes the relatively low PSNR values calculated for the images.

From the performance and quality results, we observe a benefit to compressing coverage masks prior to usage, with little visible loss in quality. The method described in Section 3 that yields these results relies heavily on 32bit integer operations but is otherwise portable to a wide variety of platforms. These performance metrics also do not take into account the possible benefits from multithreading approaches. Although these methods are highly parallelizable, the main benefit is reducing the latency of uploading the coverage masks to the GPU. Hence, any GPU compression method that would require the data uploaded prior to compression would lose this benefit. However, if the coverage information is generated on the GPU, then our method could be used to compress the mask very quickly using only a handful of low-latency integer operations.

6 CONCLUSION, LIMITATIONS, AND FUTURE WORK

In this paper we have shown that coverage masks used for rendering 2D anti-aliased non-convex paths are perfect candidates for real-time compression. Their low-entropy properties make compression algorithms very efficient and the masks themselves highly compressible. We have also shown that these masks can be compressed in real-time often speeding up the rendering of 2D curves and saving valuable GPU memory.

Limitations: Although the coverage masks can be compressed effectively, GPU-based methods for rendering arbitrary 2D-curves with anti-aliasing are still slower than their CPU-based counterparts. In general, generating the coverage mask is by far the most expensive operation of the rasterization procedure. During CPU-based rendering, the rasterizer can perform the shading directly from the RLE buffer discussed in Section 3.2. This limitation can be observed from the time it takes to run the polygon benchmark from Table 1 on different platforms using the software renderer:

Rendering time for convex path benchmark strokedrects

Platform	GPU	CPU
Moto X	$6.9\mu{ m s}$	$37.6\mu{ m s}$
Galaxy Note 10.1	$3.76\mu s$	$15.5\mu{ m s}$

Rendering time for non-convex path benchmark *polygon*

Platform	GPI	CPU	
	Uncompressed	Compressed	
Moto X	163ms	137ms	83ms
Galaxy Note 10.1	171ms	136ms	46ms

However, many of the applications that require 2D rendering operate on many more primitives than non-convex 2D curves. In the table above, the GPU-based convex path rendering operation still outperforms its CPU counterpart. For this reason, it is advantageous to use a GPU-based framebuffer. As such, our method provides benefits to the least efficient aspect of GPU-based resolution independent graphics rendering.

Additionally, as we described in Section 4, our method does not create high fidelity texture compression for general purpose grayscale images. We assume coverage masks to be



PSNR	dashed	rounded	poly	text	strokep
I JINK	35.457	41.028	35.457	37.736	53.876

Fig. 13: Detailed analysis of correctness tests within Skia most heavily affected by changes to anti-aliased non-convex path rendering. From top to bottom, the images are labeled as 'dashed', 'rounded', 'poly', 'text', and 'strokep'. We observe very few artifacts due to compression. Although the pixels along the anti-aliased edges in the rendered images do contain different pixel values contributing to the relatively low PSNR values, the detail in the edges remains. Pixels in the difference image are on if the shaded values in the corresponding original and compressed images differ. Most noticeable in 'strokep', the low entropy of the coverage masks causes pixel differences only in those along the edges of the filled paths.

highly uniform with little variation along primitive edges. If these assumptions are maintained, as we see in Figure 13, then rendering using our compression technique maintains acceptable perceptual quality.

Future Work: We have shown that coverage masks are very amenable to compression. Due to the very high fidelity of the rendered images even at the highest available compression ratios (12×12 ASTC) there is ample room for even more aggressive compression formats. Encodings that support block dimensions up to 32 or 64 may still produce nice results. The compression algorithms in Section 3.2 can be extended to support even better compression ratios, which will increase both the rendering speed and memory usage. Another direction for research is the ability to generate coverage information on the GPU itself. If such a technique existed, the compositing procedure using the coverage mask could be performed at the same time as generating the coverage information itself. However, if the

coverage mask were generated on the GPU and then used as input to a second compositing pass, compressing the GPU-generated coverage masks using this technique would incur trivial cost. Due to the random-access restrictions of compressed texture formats, they are perfect candidates for massively parallel encoding. Furthermore, to combat the original artifacts from the Blinn-Phong method, conservative rasterization may be used to cover every pixel touched by the bounding triangles [24]. Such a solution could eliminate the need for CPU-side rendering entirely. Finally, the error analysis in Section 4 opens up the possibility for additional compression algorithms that may do a better job of compressing both coverage masks and general purpose data.

7 ACKNOWLEDGEMENTS

This research is supported in part by ARO Contract W911NF-14-1-0437, Samsung, and Google. The authors would also like to thank Robert Phillips, Mike Reed, and Brian Salomon for their help and guidance in understanding and interfacing with the Skia library.

REFERENCES

- J. Barros and H. Fuchs, "Generating smooth 2d monocolor line drawings on video displays," *SIGGRAPH Comput. Graph.*, vol. 13, no. 2, pp. 260–269, Aug. 1979. [Online]. Available: http://doi.acm.org/10. 1145/965103.807454
- [2] J. M. Lane and R. a. M. Rarick, "An algorithm for filling regions on graphics display devices," ACM Trans. Graph., vol. 2, no. 3, pp. 192–196, Jul. 1983. [Online]. Available: http://doi.acm.org/10.1145/357323.357326
- [3] E. Fiume, A. Fournier, and L. Rudolph, "A parallel scan conversion algorithm with anti-aliasing for a general-purpose ultracomputer," *SIGGRAPH Comput. Graph.*, vol. 17, no. 3, pp. 141–150, Jul. 1983. [Online]. Available: http://doi.acm.org/10.1145/964967.801143
- [4] I. Google, "Skia 2d rendering library," https://sites. google.com/site/skiadocs/, 2014.
- H. Zhang, D. Manocha, T. Hudson, and K. E. Hoff, III, "Visibility culling using hierarchical occlusion maps," in *Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques*, ser. SIGGRAPH '97. New York, NY, USA: ACM Press/Addison-Wesley Publishing Co., 1997, pp. 77–88. [Online]. Available: http://dx.doi.org/10.1145/258734.258781
- [6] J. Kautz, J. Lehtinen, and T. Aila, "Hemispherical rasterization for self-shadowing of dynamic objects," in *Proceedings of the Fifteenth Eurographics Conference* on Rendering Techniques, ser. EGSR'04. Aire-la-Ville, Switzerland, Switzerland: Eurographics Association, 2004, pp. 179–184. [Online]. Available: http://dx.doi. org/10.2312/EGWR/EGSR04/179-184
- [7] I. StatCounter, "Global stats, top 5 desktop, tablet & console browsers from sept 2013 to sept 2014," http: //gs.statcounter.com, 1999-2014.
- [8] K. I. Iourcha, K. S. Nayak, and Z. Hong, "System and method for fixed-rate block-based image compression with inferred pixel values," U. S. Patent 5956431, 1999.

- [9] J. Ström and M. Pettersson, "ETC2: texture compression using invalid combinations," in *Proceedings of the* 22nd ACM SIGGRAPH/EUROGRAPHICS symposium on Graphics hardware, ser. GH '07. Eurographics Association, 2007, pp. 49–54. [Online]. Available: http://dl.acm.org/citation.cfm?id=1280094.1280102
- [10] J. Nystad, A. Lassen, A. Pomianowski, S. Ellis, and T. Olson, "Adaptive scalable texture compression," in *Proceedings of the ACM SIGGRAPH/EUROGRAPHICS conference on High Performance Graphics*, ser. HPG '12. Eurographics Association, 2012, pp. 105–114.
- [11] M. C. Shebanow, "The evolution of mobile graphics and the potential impact on interactive applications," in *Keynote Address of the ACM SIGGRAPH Symposium* on Mobile Graphics and Interactive Applications, ser. SIG-GRAPH ASIA '14. ACM, 2014.
- [12] P. Krajcevski and D. Manocha, "Compressed coverage masks for path rendering on mobile gpus," in *Proceedings of the 19th Symposium on Interactive 3D Graphics and Games.* ACM, 2015, pp. 101–108.
- [13] C. Loop and J. F. Blinn, "Resolution independent curve rendering using programmable graphics hardware," in July 2005 Transactions on Graphics (TOG) Volume 24 Issue 3 (Siggraph 2005). Association for Computing Machinery, Inc., 2005. [Online]. Available: http://research.microsoft.com/apps/pubs/ default.aspx?id=78197
- [14] M. J. Kilgard and J. Bolz, "Gpu-accelerated path rendering," ACM Trans. Graph., vol. 31, no. 6, pp. 172:1–172:10, Nov. 2012. [Online]. Available: http://doi.acm.org/10.1145/2366145.2366191
- [15] Z. Qin, "Vector graphics for real-time 3d rendering," Ph.D. dissertation, University of Waterloo, 2009.
- [16] Y. Kokojima, K. Sugita, T. Saito, and T. Takemoto, "Resolution independent rendering of deformable vector objects using graphics hardware," in ACM SIGGRAPH 2006 Sketches, ser. SIGGRAPH '06. New York, NY, USA: ACM, 2006. [Online]. Available: http://doi.acm.org/10.1145/1179849.1179997
- [17] C. Green, "Improved alpha-tested magnification for vector textures and special effects," in ACM SIGGRAPH 2007 Courses, ser. SIGGRAPH '07. New York, NY, USA: ACM, 2007, pp. 9–18. [Online]. Available: http://doi.acm.org/10.1145/1281500.1281665
- [18] G. He, B. Bai, Z. Pan, and X. Cheng, "Accelerated rendering of vector graphics on mobile devices," in *Human-Computer Interaction. Interaction Platforms* and *Techniques*, ser. Lecture Notes in Computer Science, J. A. Jacko, Ed. Springer Berlin Heidelberg, 2007, vol. 4551, pp. 298–305. [Online]. Available: http://dx.doi.org/10.1007/978-3-540-73107-8_33
- [19] A. C. Beers, M. Agrawala, and N. Chaddha, "Rendering from compressed textures," in *Proceedings* of the 23rd annual conference on Computer graphics and interactive techniques, ser. SIGGRAPH '96. ACM, 1996, pp. 373–378. [Online]. Available: http://doi.acm.org/ 10.1145/237170.237276
- [20] P. Krajcevski, A. Lake, and D. Manocha, "FasTC: accelerated fixed-rate texture encoding," in *Proceedings* of the ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games, ser. I3D '13. ACM, 2013, pp.

137–144. [Online]. Available: http://doi.acm.org/10. 1145/2448196.2448218

- [21] J. Ström and T. Akenine-Möller, "iPACKMAN: high-quality, low-complexity texture compression for mobile phones," in *Proceedings of the ACM SIGGRAPH/EUROGRAPHICS conference on Graphics hardware*, ser. HWWS '05. ACM, 2005, pp. 63– 70. [Online]. Available: http://doi.acm.org/10.1145/ 1071866.1071877
- [22] P. Wennersten and J. Ström, "Table-based alpha compression." Computer Graphics Forum, vol. 28, no. 2, pp. 687–695, 2009. [Online]. Available: http://dblp.uni-trier.de/db/journals/cgf/ cgf28.html#WennerstenS09
- [23] C. E. Shannon, "A mathematical theory of communication," *The Bell System Technical Journal*, vol. 27, pp. 379–423, 623–656, July, October 1948.
 [Online]. Available: http://cm.bell-labs.com/cm/ms/ what/shannonday/shannon1948.pdf
- [24] T. Akenine-Möller and T. Aila, "Conservative and tiled rasterization using a modified triangle setup," J. Graphics Tools, vol. 10, no. 3, pp. 1–8, 2005. [Online]. Available: http://dx.doi.org/10.1080/ 2151237X.2005.10129198



Pavel Krajcevski received the MS degree in Computer Science from the University of North Carolina at Chapel Hill and two BS degrees in Mathematics and Computer Science from the University of Chicago. Pavel is currently a PhD candidate at UNC-Chapel Hill and has previously worked full-time at Disney Interactive Studios, and as an intern at Intel, Samsung, and Google. His research interests lie in texture compression techniques and image representations that map well to graphics hardware.



Dinesh Manocha is currently Phi Delta Theta/-Matthew Mason Distinguished Professor of Computer Science at the University of North Carolina at Chapel Hill. He received his B. Tech. degree in Computer Science and Engineering from the Indian Institute of Technology, Delhi in 1987; Ph.D. in Computer Science at the University of California at Berkeley in 1992. He has coauthored more than 400 papers in the leading conferences and journals on computer graphics, robotics, and scientific computing. He has also

served as program chair for many conferences and editorial board member for more than 12 leading journals in computer graphics, robotics, geometric computing, high performance computing, and applied algebra. Some of the software systems related to collision detection, GPU-based algorithms and geometric computing developed by his group have been downloaded by more than 500,000 users and are widely used in the industry. Manocha has received awards including IBM Fellowship, Alfred P. Sloan Fellowship, NSF Career Award, Office of Naval Research Young Investigator Award, Hettleman Award at UNC Chapel Hill, and 14 best paper awards at leading conferences. He is a Fellow of ACM, AAAS, and IEEE and received Distinguished Alumni Award from Indian Institute of Technology, Delhi.