

Toward Real Time Fractal Image Compression Using Graphics Hardware

Ugo Erra

ISISLab - Dipartimento di Informatica ed Appl. "R.M. Capocelli",
Università degli Studi di Salerno, 84081 Baronissi, Italy
Email: ugoerr@dia.unisa.it

Abstract. In this paper, we present a parallel fractal image compression using the programmable graphics hardware. The main problem of fractal compression is the very high computing time needed to encode images. Our implementation exploits SIMD architecture and inherent parallelism of recently graphic boards to speed-up baseline approach of fractal encoding. The results we present are achieved on cheap and widely available graphics boards.

1 Introduction

Fractal compression is a lossy compression method introduced by Barnsley and Sloan [1] for compactly encoding images. The main idea of fractal compression is to exploit local self-similarity in images. This permits a self-referential description of image data to be yielded, which shows a typical fractal structure when decoded. This type of redundancy is not exploited by tradition image coders and this is one of the motivations for fractal image coding.

The general approach is firstly to subdivide the image using a fixed partition in simple case or adaptive partition in an advanced approach, and then to find the best matching image portion for each part. This searching phase is known to be the most time consuming part and numerous strategies has been presented to speed-up encoding. On the other hand, fractal image compression offers interesting features like fast decoding, independent-resolution and good image quality at low bit-rates which is useful for off-line applications.

Today's GPUs (graphics processing units) have high-bandwidth memories and more floating-point units. One of the most recently presented GPUs, the NVIDIA 7800, has peak performance of 165 Gflops and memory bandwidth of 38.4 GB/sec. The growing rate of the GPUs has been estimated to be 2.5 - 3.0 times a year, which is faster than Moore's Law for CPUs. Recently, all this computational power has become cheap and widely available. As side effects, several researches has began to exploit GPUs for general purpose applications such as scientific computation, database operations, matrix multiplications and many more (the growing interest in the field is witnessed by the highly dynamic website on General Purpose computation on GPUs <http://www.gpgpu.org>).

This paper presents a novel approach to perform fractal compression on programmable graphics hardware; to our knowledge, this is the first application that

uses the GPU for image compression. Using programmable capabilities of the GPUs, we exploit the large amount of inherent parallelism and memory bandwidth to perform fast pairing search between portions of the image. As a result, we show that GPUs are an effective co-processors for fractal image processor.

2 Fractal compression

The basic idea of fractal compression is to find similarities between larger and smaller portions of an image. This is accomplished partitioning the original image into blocks of fixed size, called *range* and creating a shape codebook from the original image of double size of the range, called *domain*. Range blocks partition the image so that every pixel is included while the domain blocks can be overlapped and/or to not contain every pixel.

For each range block R we must find a domain D from codebook such that $R \approx sD + o\mathbf{1}$ where s and o are called *scaling* and *offset* respectively. These values define the optimal transformation by which we can encode an image portion using another part. The encoder must scan all the codebook to find optimal D , s , and o . The domain block must be shrunk by pixel averaging to match the size of range block.

The method of least squares to find the optimal coefficients can be used. Given the two blocks R and D with n pixel intensities, r_1, \dots, r_n and d_1, \dots, d_n , the quantity to minimize is $\sum_{i=1}^n (s \cdot d_i + o - r_i)^2$ where coefficients s and o are given by

$$s = \frac{n(\sum_{i=1}^n d_i r_i) - (\sum_{i=1}^n d_i)(\sum_{i=1}^n r_i)}{n \sum_{i=1}^n d_i^2 - (\sum_{i=1}^n d_i)^2} \quad o = \frac{1}{n} \left(\sum_{i=1}^n r_i - s \sum_{i=1}^n d_i \right) \quad (1)$$

The values s, o , and the position of domain block D are the encoded values for range R . The steps of the baseline encoder with fixed block are the following:

1. *Range blocks* R_i . Given a fixed size (4×4 , 8×8 , and so on) create a set of range blocks overlapping the entire image.
2. *Shape codebook* D_i . The shape codebook is created in two steps:
 - (a) Using a step size of l pixel horizontally and vertically create a set of domain blocks which are double the range size.
 - (b) Shrink the domain blocks by averaging four pixel to match range size.
3. *The search*. For each range block R an optimal approximation $R \approx sD + o\mathbf{1}$ is computed in the following steps:
 - (a) For each domain block D_i compute $R \approx sD_i + o\mathbf{1}$ performing the least square optimization formula using formulas (1).
 - (b) Among all codebook D_i output the code for current range $[k, s, o]$ such that the error $R \approx sD_k + o\mathbf{1}$ is minimum.

The output encoder is a description of an operator which serves as approximation of the original image. An operator T is defined over an image f as

$Tf \equiv sf + o\mathbf{1}$. Thus, starting from any initial image g_0 and applying interactively T to obtain $g_1 = Tg_0$, $g_2 = Tg_1$, $g_3 = Tg_2, \dots$ the sequence g_i converges to an approximation g_n of the original image called *attractor* after few iterations. The mathematical theory about these principles can be found on [2].

Related works. Fractal compression allows fast decompression but has long encoding times. The most time consuming part is the domain blocks searching from each range.

In order to reduce the number of computations required, some authors have proposed to restrict the set of domain blocks. For instance, in [3] Beaumont adopts a search strategy using an outward spiral starting from the coordinate of range and halts when a necessary condition has been reached. These strategies reduce encoding time but image quality could suffer due to the overlook of some possible optimal pairing.

Categorized search proposed by Boss, Fisher and Jacobs [4, 5] and features vector methods proposed by Saupe [6] are efficient classification techniques. They reduce the encoding complexity using a classification of the domain codebook block in such way that for each range the search is essentially more efficient.

The use of general purpose high performance architecture has been used to accelerate the encoding phase without a decrease of image quality.

Related work has been done concerning the encoding phase on SIMD architecture. In [7] massively parallel processing approach has been used on an APE100/Quadrics SIMD machine. For testing, they used 512 floating point processors, offering a peak power of 25.6 GFLOPS. They are able to compress a gray level image of 512x512 using a scalar quantization techniques in about 2 seconds. A pixel-based parallelization scheme as been used in [8] on a pyramidal SIMD architecture for a non-adaptive version of fractal compression. While, in [9] an image-block parallelization scheme is used on a SIMD array processor.

3 Programmable graphics hardware

Today, graphics cards are fundamentally programmable stream processors [10]. In this computational model stream are collections of data requiring similar computation. Every object in the stream is processed by the some function called kernel. The stream elements that GPUs process are vertices and fragments, which are essentially pixels with some related information. Instead, the kernels that GPUs are capable to execute are called vertex and fragment programs. Fragment processor is designed to run fragment programs and expose features that are more suitable for general purpose application.

The fragment processor computes pixel color as an RGBA data. Notice that the set of four-wide SIMD floating-point operations is sufficiently complete and flexible. Furthermore, as stream architecture, the fragment processor exploits spatial parallelism; it runs the same fragment program for each pixel. In the case of the last graphic cards, up to twenty four-pixel pipelines handle single-instruction, multiple-computing processing.

Fragment processor is the last stage of graphics pipeline. Each fragment, after processing, ends up as an RGBA pixel on a screen. Alternatively, using a p-buffer, the pixels can be captured in off-line image buffer as a texture. This model presents limits and advantages. For each incoming pixel the fragment program is invoked at a specific location and returns the final value in the same location as output. That is not possible to write in a different location or to do scattering. Instead the kernel can do gathering using textures as lookup table to read precomputed values.

The textures are an invaluable resource because they can be used as lookup tables during computation though access to them is restricted to write-only or read-only. Floating-point texture vectors can be of two, three or four components. Each fragment can fetch a component vector as input from one or more textures and returns vector components as output up to four textures. Today, it is possible to find graphics boards that support 16-bit, 24-bit, and 32-bit floating point per color components. Thus, we could fetch an unlimited 256 bit values and return up to four 256 bit values.

4 Mapping fractal compression on the GPU

In order to exploit the specialized nature of the GPU and its restrict programming model we must map the fractal compression as a streaming computation. The goal is to perform pairings test between range and domain exploiting parallel architecture of the GPU and high bandwidth access to pixels. The entire process uses a gray level image as input data and returns the textures T_{POS} with the position of optimal domain blocks and T_{SO} with scaling/offset coefficients as outputs.

The underlying idea is to use a producer/consumer scheme. The producer gathers from the domain pool a block which is broadcast to all consumers that are the ranges. Each range stores the current domain as soon as it appears as the best pairing block. The entire process continues until all domain blocks have been consumed. In this scenario, a pixel appears as a single floating-point processor responsible for only one range. Then, the GPU mimics a computational grid, as shows in Figure 1, rendering a sized-range rectangular upon which performs parallel pairing test among all the ranges for a given domain.

The next two subsections describe the organization of data structures on the GPU and the fragment programs to accomplishment the entire process.

4.1 Data Structures Organization

Fractal compression is implemented as fragments programs. These programs are executed via multi-pass rendering of a screen-sized rectangle where each pixel is an encoding range. Notice that during encoding the same range will be paired among all other domains and from another point of view the same domain will be paired among all ranges. Then, for each range block and for each domain block we precompute all the related quantities that remain constants during the

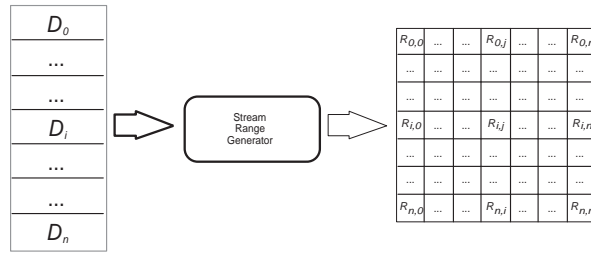


Fig. 1. A computational grid based on GPU. For each rendering pass a domain block is broadcast to all ranges.

entire encoding. These constant values are the summations of the scaling and offset formulas in 1 and will be stored in the lookup textures T_R and T_D as shows in Table 1.

In order to exploit SIMD parallelism and efficient bandwidth, during the encoding, the source image is stored compactly into texture. An RGBA texture which uses 32-bit per component is capable to store up to 256 bit per pixel. Usually, for a gray level image are necessary 8-bit per pixel. Then, thanks specialized instruction `pack` we are able to store up to 16 pixels into a single RGBA pixel's texture. Using this representation of the image it is possible to read 16 pixels simultaneously in a single texture access followed by a `unpack` instruction.

Table 1. Look-up textures for constant summations. In texture T_D we store two summations into components red and green

Summations	Texture	Component Used
$\sum_{i=1}^n r_i$	T_R	R32
$\sum_{i=1}^n d_i$	T_D	R16
$\sum_{i=1}^n d_i^2 - (\sum_{i=1}^n d_i)^2$	T_D	G16

4.2 Fractal compression kernels

The entire flow diagram for the streaming fractal compression is illustrated in Figure 2. The following sections detail the implementation of each kernel and the read/write textures access. All the following kernels always take as input the compact version of source image.

Ranges summation. This kernel precomputes the range summation using the technique described in the previous section. It takes the original image as input and returns a texture T_R as output. This buffer is a previously declared one 32-bit color component texture. The size of this buffer is 1/4 of the input

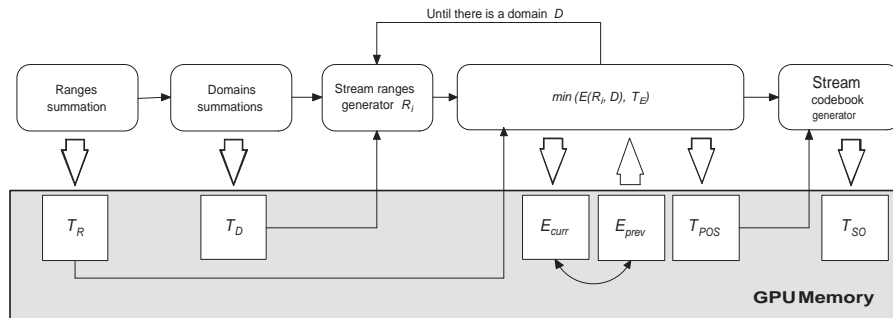


Fig. 2. A streaming fractal compressor.

image if we choose a range block size of 4×4 , or $1/8$ for a range block of 8×8 , and so on.

Domains summation. The kernel precomputes the domains summations. It takes the original image as input and returns the texture T_D as output. This buffer is a previously declared two 16-bit color components texture.

Stream range generator. This is the only party performed in the CPU, it is not a computational stage but it serves as a start-up routine to generate a stream of fragment programs. It draws a texture of size T_R to force the “Pairing Test” fragment program execution for each range. Furthermore, it is responsible to pass the position of current domain and the entire set of pixels belongs to the current domain as parameters. This permits to store an entire domain into registers of fragment processor avoiding continuous texture fetches of the same pixels’ domain. The stream range generator is invoked in a multipass rendering for each domain block.

Pairing test. This kernel performs all possible tests for optimal pairing. At each rendering pass this kernel has T_R as input textures, coordinate of current domain and pixels’ domain from the previous stage. It returns a sized range texture T_{POS} with the coordinates of current optimal blocks as output. More precisely, due to hardware limits which not allow read/write access on the same texture, pairing tests are not performed for a given range R among all domains D_i . We must invert pairing tests as follows: given a domain D performs pairing tests among all ranges R_i updating minimal errors for each range. This schema use a double buffer to read previous errors in texture T_{prev} and to store current errors in texture T_{curr} . Thus, given a domain block D at each rendering pass each fragment program computes and stores the following value: $T_{curr} = \min(E(R_i, D), T_{prev})$. Before the next rendering pass the errors textures T_{curr} and T_{prev} are swapped and another domain block is passed as input.

Stream codebook generator. The last stage computes and writes into texture T_{SO} final scaling s and offset o . The kernel has textures T_{POS} as input and returns a texture T_{SO} with optimal coefficients. This operation is performed here

in order to avoid useless write operations and therefore optimizes the amount of memory bandwidth required for texture accesses.

5 Fractal decomposition

The decoding phase is very simple and is performed using only one kernel. It performs the decoding via n multipass rendering from any initial image g_0 . The fragment encoder takes the texture of optimal position T_{POS} , scaling/offset texture T_{SO} and an intermediate texture image g_i as input and returns the image g_{i+1} . The kernel works in a feedback fashion, it takes g_i and returns g_{i+1} as output image which serve as input in the next interaction until it converges to an attractor g_n . More precisely, we perform an original image-sized rendering to decode each pixel of image g_{i+1} by kernel. For each one we fetch scaling factors s and o from T_{SO} and its code block position from T_{POS} which is used to fetch pixel from g_i . In the end, the entire process takes only four texture accesses.

6 Experimental Results

In order to compare the amount of pairing tests that GPU is capable to perform, we implemented heavy brute force algorithm on GPU as well as on CPU. We have experimented on a Pentium IV based machine with 3.2GHz processor speed and 1GB of RAM. The graphics card used was a GeForce FX 6800, which has 16 internal pixel pipelines, with 128MB of video memory, core speed of 300MHz and memory speed of 800MHz. We used OpenGL to implement our GPU-based compressor.

The test image has a resolution of 256×256 pixels. Choosing a range size of 4×4 pixels we obtain 64×64 ranges. The domain blocks must be twice the size of range blocks and using a step of one pixel to scan the image, the domain pool contains $(128 - 4 + 1)^2 = 15,625$ elements. In total, using a heavy brute force strategy $4,096 \times 15,625 = 64,000,000$ possible pairings require testing.

The CPU version takes about 280 seconds to perform all pairing test whereas the GPU version takes about 1 second. Then, the amount of paring test that the GPU is capable to perform is about 21 millions per second whereas the CPU performs about 220 thousands paring test per second. These results come from considering fractal compression a random-memory-access intensive problem. The memory bandwidth together arithmetic intensity advantages GPU over CPU. Moreover, our work shows its advantages when compared with expensive parallel architecture as for instance in [7] which use 512 floating-point processors with performance comparable to our GPU implementation.

Finally, performances of our implementation are comparable with the results obtained using a different approach as in [2], where compression time is traded-off with compression ratio (and, therefore, image quality) by using a parametric adaptive partitions. In fact, on the same hardware, the implementation in [2] compresses image approximately in same time but with lower image quality. Of course, it should be emphasized that, given the quicker growing rate of GPUs vs

CPUs, the gap between our solution and the results described in [2] is going to widen in the near future, as quicker GPUs will be available on the ordinary PC.

7 Conclusions and further work

Fractal image compression is well suited for parallel system due to its high computation complexity and regular algorithmic structure. Today, we think that graphics board offers substantial computational power to take full advantages of the baseline approach. We believe that are possible two scenarios and which we are going to investigate. The first is to further exploit the graphics hardware to obtain more efficient compression schema and taking into account also color image. The second is to use GPU as an efficient co-processor. The general purpose architecture of the CPU permits to excel on arranging data. Then, the GPU could be used as an efficient pairing engine whiles the CPU arranges pairings tests.

Today, the consumers for less than \$500 can buy an off-the-shelf graphics card with performances comparable to SIMD parallel machines that cost hundreds of thousands of dollars several years ago. Furthermore, as the GPU consumes less power than a high-end CPU, it is evident how using the graphics card can extend the life-time of an existing computer system. In conclusion, we think that GPUs offer the great opportunity to take full advantages of the fractal compression on consumer desktop personal computers.

References

1. Barnsley, M.F., Sloan, A.: Chaotic compression. *Computer Graphics World* (1987)
2. Yuval, F.: *Fractal Image Compression - Theory and Application*. Springer-Verlag, New York (1994)
3. Beaumont, J.M.: Image data compression using fractal techniques. *British Telecom Technol. Journal* **9** (1991) 93–109
4. Jacobs, E.W., Fisher, Y., Boss, R.D.: Image compression: A study of the iterated transform method. *Signal Processing* **29** (1992) 251–263
5. Yuval, F.: Fractal image compression. *Fractals: Complex Geometry, Patterns, and Scaling in Nature and Society* **2** (1994) 347–361
6. Saupe, D.: Accelerating fractal image compression by multi-dimensional nearest neighbor search. In Storer, J.A., Cohn, M., eds.: *Proceedings DCC'95 Data Compression Conference*, IEEE Computer Society Press (1995)
7. Palazzari, P., Coli, M., Lulli, G.: Massively parallel processing approach to fractal image compression with near-optimal coefficient quantization. *J. Syst. Archit.* **45** (1999) 765–779
8. Xue, M., Hanson, T., Merigot, A.: A massively parallel implementation of fractal image compression. In: *Proc. ICIP-94 IEEE International Conference on Image Processing*, Austin, Texas (1994)
9. Giordano, S., Pagano, M., Russo, F., Sparano, D.: A novel multiscale fractal image coding algorithm based on simd parallel hardware. In: *Proceedings of the International Picture Coding Symposium PCS'96*, Melbourne (1996)
10. Venkatasubramanian, S.: The graphics card as a stream computer. In: *SIGMOD-DIMACS Workshop on Management and Processing of Data Streams*. (2003)