

An Approach on Hardware Design for Computationally Intensive Image Processing Applications based on Light Field Refocusing Algorithm

Jiayuan Meng, Dee A.B. Weikle, Greg Humphreys, Kevin Skadron
Dept. of Computer Science, University of Virginia
{jm6dg, dweikle, humper, skadron}@cs.virginia.edu
UVA Technical Report CS-2007-15

Abstract

This paper describes the performance analysis of the light field refocusing algorithm running on different hardware specifications, including the Intel Pentium 4, SSE2(Streaming SIMD Extensions), GPU, and also Cell Broadband Engine. The hardware chosen has unique features, making it interesting to compare their performance on such an application with each other, and how much advantage or disadvantage each one has over others. By doing so, we attempt to clarify the pros and cons of each hardware design in its capability of handling computationally intensive image processing applications. The execution time is used as the main metric.

1 Introduction

Light Field Photography is a new technology which trades spatial resolution for directional resolution. It is able to capture an image which, after computation, reveals the object from different view angles and also at different depth.

With the raw image captured by the camera, the refocusing algorithm [2] computes a new image focused at a user-specified focal depth. For a 4096×4096 light field image with 256×256 spatial resolution and 16×16 directional resolution, computation can barely reach real-time with traditional spatial algorithms [2]. Moreover, the size of the image will grow in the future. In fact, many computationally intensive image processing applications are calling for realtime hardware design.

Most of these image processing applications involve a large amount of memory accesses, and the operations are often repetitive. We chose the light field refocusing algorithm as a typical and booming image processing application. The repetitive behavior denotes a high potential for parallelism. Meanwhile, a variety of hardware has been designed for parallelism, including SSE2, GPU, and multi-core chips. Although not originally designed for image processing applications, they each contribute some inspiration for future designs.

The rest part of this paper is organized as follows: Section 2 gives an overview of the light field refocusing algorithm. Section 3 presents how we implemented the refocusing algorithm on different hardware. Performance is compared and analyzed in Section 4. Conclusions are drawn in Section 5.

2 Refocusing Algorithm of Light Field Photography

2.1 A Basic Algorithm

The concept of light field photography was originally developed by Ren Ng and Marc Levoy [3]. The image captured can be regarded as a 4D light field array, with 2 dimensions for spatial resolution and another 2 for directional resolution. For a light field image of size 4096×4096 captured by a plenoptic camera with a 256×256 microlens array, the spatial resolution is 256×256 and the directional resolution is 16×16 . Behind each microlens, a small 16×16 version of the original image is formed, which actually correspond to 16×16 sub-apertures. If we extract those pixels corresponding to the same sub-aperture, then we get a sub-aperture image. In our example, the sub-aperture image is 256×256 , and there are 16×16 sub-aperture images. Figure 1 shows a sample light field image [2].

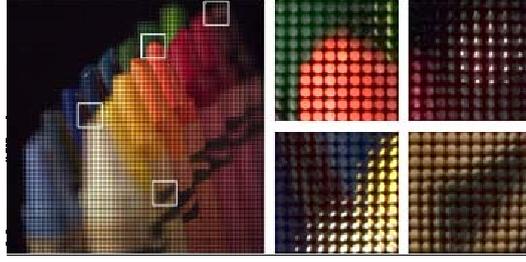


Figure 1: A sample light field image

The basic task of the refocusing algorithm is illustrated in Figure 2. Rays focusing at certain depth will be scattered and captured with different angles and different microlenses, and thus will be mapped onto a set of pixels. The refocusing algorithm computes an inverse mapping and adds up all the pixels corresponding to the same position on the refocused plane. The refocusing algorithm is described in detail in [2].

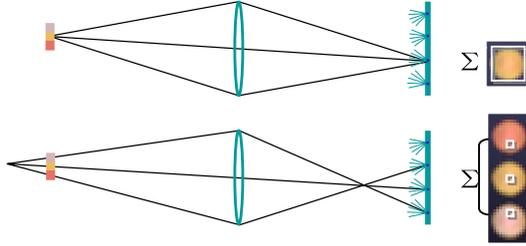


Figure 2: Refocusing Principle: lights focused at certain depth is scattered and captured by different microlenses. All we need to do is gather them and summing up!

This refocusing process is referred to by a **Photography Operator**: For a position (x, y) on the refocused plane, integration is computed for rays traveling from all the positions (u, v) on the main lens plane to (x, y) on the refocused plane. It is illustrated in Figure 3.

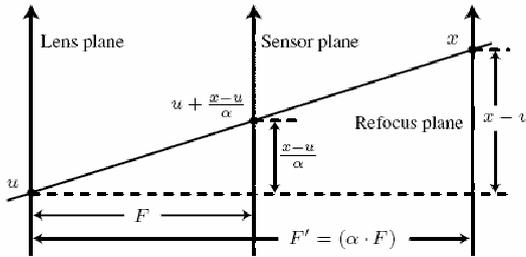


Figure 3: Illustration of the refocusing algorithm

The Photography Operator can be expressed by:

$$P\alpha = \frac{1}{\alpha^2 F^2} \int \int L_F(u(1 - \frac{1}{\alpha}) + \frac{x}{\alpha}, v(1 - \frac{1}{\alpha}) + \frac{y}{\alpha}, u, v) du dv \quad (1)$$

Where L_F is the 4D light field captured by the camera, x, y are defined on the refocused plane, u, v are defined on the main lens' plane. F is main lens' focal length. α is a scaling factor which determines the refocusing depth. $\alpha \cdot F$ is the distance between the refocused plane and the main lens' principle point. A ray is defined with a vector (x, y, u, v) . This operator then computes which position on the focal plane that the ray corresponds to, since it's the focal plane where we have a image formed. We then extract color information from the light field image captured at the focal plane.

2.2 A Derived Algorithm

Assume the spatial resolution is $N * N$. After computation with the basic algorithm, the refocused image has the size of $\alpha N * \alpha N$. This is physically correct. Imagine moving the refocus plane in Figure 3 forward and backward, the area which forms image will vary by a factor of α . If we sample the refocused image at the same rate for different refocusing depth, the effective image size will vary according to α . However, in [2], Ren showed that the Nyquist resolution is $N * N$ with exact refocusing and less with inexact refocusing. Thus, we may want to adaptively vary the sampling rate according to α to achieve the best possible resolution.

The new Photography Operator then becomes:

$$P\alpha = \frac{1}{\alpha^2 F^2} \int \int L_F(u(\alpha - 1) + x, v(\alpha - 1) + y, u, v) dudv \quad (2)$$

This Photography Operator generate refocused image with a fixed size of $N * N$.

Note that a sub-aperture image is defined as

$$S_{(u_s, v_s)}(x, y) = L_F(x, y, u_s, v_s) \quad (3)$$

This is very similar to the inner part of Equation 2. The only difference is that in Equation 2, x has an offset of $u(\alpha - 1)$, and y has an offset of $v(\alpha - 1)$. Thus, we can treat the derived refocusing algorithm as the summing of all the sub-aperture images with different offsets according to α .

Intuitively, we can think of sub-aperture images as a group of images taken from slightly different angles. Thus, objects at different depth will be shifted with different amount. The refocusing algorithm computes how much offset each sub-aperture image has for a certain refocusing depth, and then align the sub-aperture images so that objects at the refocusing depth overlap in all the images and appear focused, while objects elsewhere do not overlap and appear blurry.

The derived algorithm also improves the performance of the refocusing algorithm, since the offset is fixed for a sub-aperture image, we don't have to compute them for every element. It also exploit locality in memory access. Instead of fetching data ad-hoc, we can process the sub-aperture images in sequence.

This has a tremendous improvement in performance (see Table 1 in Section 4) Since we are targeting at an optimized algorithm as our base for further optimization on different hardware, the derived algorithm is chosen for further development.

Moreover, we used only nearest point interpolation. Note that the offset is a floating point, bilinear interpolation or more complex interpolation scheme can be applied. Figure 4 compares the image quality between nearest point interpolation and bilinear interpolation. The raw image is from light field microscopy[1] The results doesn't show much difference. Since complex interpolation significantly elongate the execution time(see Table 1 in Section 4), we simply used nearest point interpolation, which just round the floating point to the nearest integer.

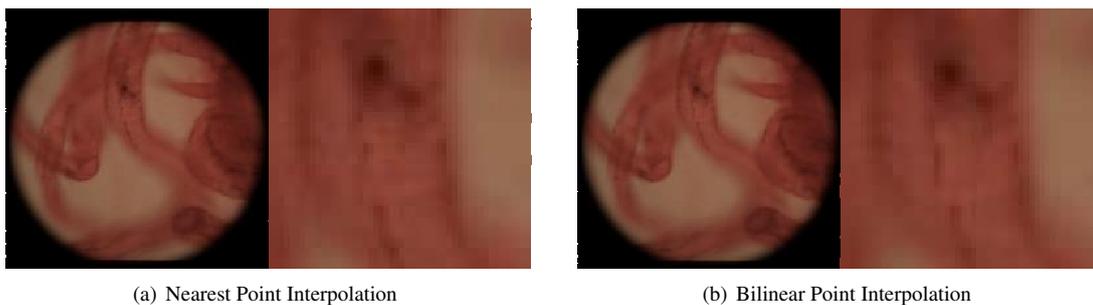


Figure 4: Comparison of Nearest Point Interpolation and Bilinear Point Interpolation

A more detailed comparison of the implementation is shown in Figure 5.

3 Implementation on Different Hardware

We implemented the refocusing algorithm with regular C code, SSE2 instructions, GPU fragment shader, and Cell Broadband Engine. A brief introduction of GPU and CBE can be found in [4] and [5]. The plain C code has been shown as the

Derived Refocusing Algorithm in Section 2, and is run on Intel Pentium 4 as a standard reference. Let's now focus on the other 3 implementations.

3.1 Implementation with SSE2

SSE2 is specialized in vectorized stream data. It has multiple 128 bit registers. In our implementation, data is added column by column. Each line is processed with SSE2 instructions. Since the offset varies for different sub-aperture image, it is not possible to align the data once and for all. Therefore, we use those instructions for unaligned vectors. The code and annotation is shown in Figure 6

3.2 Implementation with Cell Broadband Engine

With CBE [5], we use the master-slave programming model. The refocused image is divided into 8 portion, each portion is computed by an SPU. There are two pieces of codes, one for PPU, which divides the work and synthesize it at last. The main code for refocusing algorithm is on SPU. Since CBE processes data with 128 bit aligned vectors, we have to shuffle the unaligned data. To ease the task, we reorganize each color data as 4 shorts. 3 of them representing RGB channels and the last is used just for padding. A short in SPU is 16 bits long, the same as char16. With this modification, a 128bit vector contains 2 pixels' information. We are using the CBE simulator for testing the performance. We count the cycles and then convert it to time according to the clock rate. Code is provided in Figure 7 and Figure 8

3.3 Implementation with GPU Fragment Shader

We use CG[4] to program the fragment shader on GPU. The performance of GPU depends heavily on the hardware. Since the CG compiler has to unroll the loops, our fragment shader is limited by not only texture size, but also the code size. This limitation depends on specific GPUs. For GPUs with small capacity, we have to do multipass, which increases complexity and has a significant overhead in performance. We will demonstrate the performance in Section 4.

As a sample code, we demonstrate an implementation with 64 passes, in each pass, $2 * 2$ sub-aperture images are processed. Implementation with less passes only requires a few modifications on the global constants.

The initialization prior to the refocusing is straightforward. The program reads in the light field image and extract sub-aperture images as before. The sub-apertures images are then tiled into groups to form textures for GPU processing. In our example, a texture is a square image composed of $2 * 2$ sub-aperture images.

We also used render-to-framebuffer technique to store the result from former passes. Figure 9 describes the code in detail.

4 Performance and Analysis

Table 1 compares the performance of the Basic Algorithm and the Derived Algorithm. We also compare their performance with gcc's sse2 compiler option[-msse2]. In addition, bilinear interpolation is also compared.

<i>Spatial&Directional</i>	<i>B</i>	<i>Bs</i>	<i>D</i>	<i>D(-msse2)</i>	<i>D(bi)</i>	<i>D(bi,-msse2)</i>
(200*200),(20*20)	1.219	1.047	0.266	0.266	3.36	4.422
(256*256),(16*16)	1.187	0.984	0.297	0.281	3.562	4.672

Table 1: Performance Comparison of the Basic Algorithm(B) and Derived Algorithm(D), as well as the Derived Algorithm with bilinear interpolation(bi). Execution time in seconds are used as metrics.

We see that the gcc's [-msse2] compiler option does improve the performance for the Basic Algorithm. However, it hardly improves the performance of the optimized Derived Algorithm. It even produces significant overheads in some cases, e.g. when doing bilinear interpolation.

Indeed, using SSE2 is not always beneficial. SSE2 is based on vectors for stream data. However, many image processing applications are not stream data. Many of the applications require windowing and involves a large amount of computation in locality. Implementing this kind of operation in SSE2 requires additional complex operation. In many case a vector has

to be split or permuted before computation. SSE3 and SSE4 have added instructions to add and subtract the multiple values stored within a single register.

Table 2 compares the performance of the algorithm with different hardware specifications described in Section 3. Since GPU can only hold textures with a size of a power of 2, we have to pad the image with spatial resolution (200*200) to (256*256). The directional resolution remain the same.

<i>Spatial&Directional</i>	<i>SSE2</i>	<i>GPU_{64passes}</i>	<i>GPU_{16passes}</i>	<i>GPU_{4passes}</i>	<i>CBE</i>
(200*200),(20*20)	0.047	0.026	0.026	0.026	0.011
(256*256),(16*16)	0.047	0.013	0.013	0.026	0.00434

Table 2: Performance Comparison of the Basic Algorithm(B) and Derived Algorithm(D), as well as the Derived Algorithm with bilinear interpolation(bi). Execution time in seconds are used as metrics.

Compare the performance of the manually tuned SSE2 instructions with the compiler generated SSE2 instructions in Table 1, we can see that the manually tuned SSE2 instructions can still improve the performance significantly.

Cell Broadband Engine is a multi-core exploitation of parallelism. Similar to SSE2, it is also based on vectors. The 8 SPU's compute data in parallel, so it is even faster than SSE2, and is comparable with GPU. However, with 8 SPU's, the performance is not 8 times that of SSE2. This difference may due to branch miss penalty as well as the Direct Memory Access overhead. CBE tries to alleviate these negative impacts by introducing branch hint and double buffering. To process unaligned data, we have to shuffle the bits in different vectors. However, this mechanism also requires many branches, which introduce a branch penalty which is hard to be alleviated by branch hint.

GPU provides options for a multitude of vector sizes. It is also optimized for texture access in locality. Moreover, it supports interpolation in hardware level, so the GPU implementation is automatically interpolated. However, in general, GPU is not as flexible as CBE in many aspects. First, GPU has a fixed programming model — the graphics pipeline. Parallelism can only be exploited within certain stages, and the resource is not fully utilized for many GPGPU applications. In our example, only fragment processors are utilized. Moreover, in a single pass, memory writes is limited to vertex program, and memory reads is limited to fragment program. Multipass is usually implemented for more complex memory access demands. In addition, GPU performance depends heavily on specific graphics card. Code has to be tuned for capacity of the particular card with respect of texture size and instruction size. Our experiment is based on NVidia GeForce 7600GS model. We have tried to tile the sub-aperture images to textures with sizes of 512*512, 1024*1024 and 2048*2048, and for each size of texture, we need to do 64, 16, and 4 passes. According to Table 2, the performance of 64 passes and 16 passes are identical. However, with a texture of 2048*2048, the performance drops even if only 4 passes are needed. This is due to the overhead of image transferring between the main memory to the memory on GPU.

5 Conclusion and Future Work

With respect to image processing applications, we list the strengths and weaknesses of SSE2, CBE and GPU as follows:
SSE2:

- *Strengths:*
 1. *Vector level parallelism, work best for stream data*
 2. *Simplicity, integrated well into the single chip, thus less memory transfer and branch miss penalty*
- *Weaknesses:*
 1. *No chip-level parallelism*
 2. *Difficult to process small vectors(like those with 32 bits or less)*

CBE:

- *Strengths:*
 1. *Exploited parallelism in both vector level and chip level*

2. Enable the user to control memory transfer as well as branch hint, give more opportunity to optimize applications
3. Provide flexible programming model for more complex applications. In addition, user can utilize the resources as much as possible

- Weaknesses:

1. Vectors are most suitable for stream data, operating smaller vectors requires much more effort. Flexibility has increased complexity.
2. Although efforts has been made to accelerate memory access, the communication between the PPU and SPUs still accounts for a long time.

GPU:

- Strengths:

1. Exploit parallelism both in multi-chips and in vectors with 16,24,and 32 bits.
2. Locality in memory is handled nicely
3. More basic functions are hardware supported, such as sin/cos, and bilinear interpolation.
4. accessorial processors are closely integrated into the pipeline, there is less memory transfer.

- Weaknesses:

1. Pipeline is fixed, resources cannot be fully utilized. Additional efforts such as multi-pass is required when more complex computation is involved.
2. Limited memory access in different processing units

Thus, we conclude that an ideal hardware model should have the following features:

- *Scalable Parallelism*: Vector size varies from 16 bit to 128 bit or more. Chip-level parallelism also benefits given a good integration
- *Simplicity with Speciality*: With simplicity comes improvement in speed. A simple unit can achieve the best performance for its specialized task. However, simple units alone has limited capability, as described in SSE2 analysis.
- *Flexibility*: Flexible combination of operations provide possibilities for more complex applications. However, with flexibility comes complexity. As for the example of CBE, the SPU provide a variety of functions, however, this also slows down its performance. We often have to trade off simplicity with flexibility.
- *Locality-aware Memory Units*: A special memory units which can tile a texture into blocks for the ease of computation in locality.

In the future, a possible method to achieve simplicity, speciality and flexibility at the same time may be that of a pool of simple but specialized accessorial functional units with programmable connections. The units can be simple vector registers and ALUs as in SSE2, or can be small processors as vertex and fragment processors in GPU. Programmable connection between the units can result in better resource utilization as well as flexible pipeline for complex applications. Small units also alleviate the burden of memory transferring.

6 Appendix: Code Annotation

References

- [1] M. Levoy, R. Ng, A. Adams, M. Footer, and M. Horowitz. Light field microscopy. *ACM SIGGRAPH*, 2006.
- [2] R. Ng. Fourier slice photography. *ACM SIGGRAPH*, 2005.

- [3] R. Ng, M. Levoy, M. Brdif, G. Duval, M. Horowitz, and P. Hanrahan. Light field photography with a hand-held plenoptic camera. *Stanford University Computer Science Tech Report CSTR 2005-02*.
- [4] J. Owens and D. Luebke. A survey of general-purpose computation on graphics hardware. 2005.
- [5] I. Systems and T. Group. Cell broad band engine programming tutorial. Oct. 2005.

<pre> float beta = (1.0f-1.0/alpha)*alpha; int weight = nu*nv; for(int x=0; x<nx; ++x) for(int y=0; y<ny; ++y){ r=g=b=0; for(u=0; u<nu; ++u){ kux = u*beta; for(v=0; v<nv; ++v){ kvy = v*beta; kx = (int) (x+kux); ky = (int) (y+kvy); ku = (int) u; kv = (int) v; if(kx<0 ky<0 kx>=nx ky>=ny) continue; If_pixel=GetColor(&lightfield, kx, ky, ku, kv); r += (unsigned int) If_pixel[0]; g += (unsigned int) If_pixel[1]; b += (unsigned int) If_pixel[2]; } } r = (float)r/weight; g = (float)g/weight; b = (float)b/weight; CLAMP(r); CLAMP(g); CLAMP(b); pixel = GetPixel(&refoc_img, x, y); pixel[0] = (unsigned char) r; pixel[1] = (unsigned char) g; pixel[2] = (unsigned char) b; } } </pre>	<p>coefficient in refocusing equation value used for normalizing the color for each pixel on refocused image r, g, b: unsigned int for color components for each sub-aperture kux, kvv: float temporaries compute where on the focal plane does the ray hits (equation 2) kx, ky, ku, kv: indices to the 4D lightfield make sure it is in a valid range get color from lightfield's 4D array If_pixel: a pointer to unsigned char add them up normalize the color clamp to 0-255 assign the color to the pixel in the refocused image</p>
--	---

(a) Code of Basic Refocusing Algorithm

<pre> beta = (1.0f-1.0/alpha)*alpha; for(u=0; u<nu; ++u) for(v=0; v<nv; ++v){ offset_x = ((float)u)*beta; offset_y = ((float)v)*beta; beginx = (offset_x<0)?-offset_x : 0; endx = (offset_x<0)?nx - nx-offset_x; beginy = (offset_y<0)?-offset_y : 0; endy = (offset_y<0)?ny - ny-offset_y; for(x=beginx; x<endx; ++x){ src = &sub_images[3*(beginy+offset_y+ ny*(x+offset_x+nx*(v+nv*u)))] ; dst = &refocus_dst[3*(beginy+ny*x)]; for(y=beginy; y<endy; ++y){ dst[0]+=(float)src[0]; dst[1]+=(float)src[1]; dst[2]+=(float)src[2]; dst+=3; src+=3; } } } } JSAMPLE* pixel; * If_pixel; float w = F/nu/nv; for(x=0; x<nx; ++x){ for(y=0; y<ny; ++y){ dst = &refocus_dst[3*(y+ny*x)]; src = GetPixel(&refoc_img, x, y); src[0] = (JSAMPLE)CLAMP(dst[0]*w); src[1] = (JSAMPLE)CLAMP(dst[1]*w); src[2] = (JSAMPLE)CLAMP(dst[2]*w); } } } </pre>	<p>Coefficient in refocusing equation For each sub-aperture image Compute the offsets Compute where to start and end for each line and row For each pixel on the sub-aperture image sub_images: unsigned char*, the merged sub-aperture images src: unsigned char* dst: float* Add their value to the corresponding position on the refocused image ∞ Since each pixel on the refocused image is added multiple times, JSAMPLE: unsigned char w is used to normalize the value Clamp the value between 0 and 255</p>
--	--

(b) Code of Derived Refocusing Algorithm

Figure 5: Comparison of the Basic Refocusing Algorithm and the Derived Refocusing Algorithm

<pre> int offset_x, offset_y; char16* dst; char16* src; char16* sub_image; int beginx, endx, beginy, endy, startrow, endline; int step_row = ny*3; for(int u=0; u<nu; ++u){ for(int v=0; v<nv; ++v){ sub_image = \ &sub_images[3*(ny*nx*(v+nv*u))]; offset_x = ((float)u)*beta; offset_y = ((float)v)*beta; begin = (offset_x<0)?-offset_x : 0; endx = (offset_x<0)?nx : nx-offset_x; beginy = (offset_y<0)?-offset_y : 0; endy = (offset_y<0)?ny : ny-offset_y; startrow = (offset_y<0)?0 : offset_y; endline = (endy-beginy)*3; src = &sub_image \ [3*(startrow+ny*(beginx+offset_x))]; dst = \ &refocus_dst[3*(beginy+ny*beginx)]; for(x=beginx; x<endx; ++x){ int i; int divider=8; asm("intel_syntax noprefix\n"); asm(" mov ecx, %0\n"; : "r"(divider) : "%ecx"); asm(" mov edx, %0\n"; : : "%edx", "%ecx"); asm(" mov eax, %0\n"; : "r"(endline) : "%eax", \ "%ecx", "%edx"); asm(" div ecx\n"); asm(" mov ebx, eax\n"); asm(" mov ebx, %0\n"; : "r"(dst) : "%eax", \ "%ebx", "%ecx", "%edx"); asm(" mov eax, %0\n"; : "r"(src) : "%eax", \ "%ebx", "%ecx", "%edx"); </pre>	<p>Color components are defined to be 16 bits, which is enough for computing colors based on RGB values in [0,255], and can utilize the vector instruction as much as possible</p> <p>For each sub-aperture image</p> <p>sub_images is a 1D array representing the sub-aperture images extracted from the 4D lightfield with a dimension of (nu, nv, nx, ny). Each element in the array is a triple char16 consisting RGB values. To index data at position (x,y) in sub-aperture images(u, v), the indexing method is sub_images[3*(y+ny*(x+nx*(v+nv*u)))]</p> <p>Compute the offsets, as well as where to begin and end for each line and row</p>
<pre> asm(" loop\n"); asm(" movdqu xmm0, [eax]\n"); asm(" movdqu xmm1, [ebx]\n"); asm(" padd xmm0, xmm1\n"); asm(" movdqu [ebx], xmm0\n"); asm(" add eax, 16\n"); asm(" add ebx, 16\n"); asm(" sub ecx, 1\n"); asm(" jnz loop\n"); asm(" mov ecx, edx\n"); asm(" shr ecx\n"); asm(" leftover\n"); asm(" jz endloop\n"); asm(" mov edx, [ebx]\n"); asm(" add edx, [eax]\n"); asm(" mov [ebx], edx\n"); asm(" add ebx, 4\n"); asm(" add ebx, 4\n"); asm(" sub ecx, 1\n"); asm(" jmp leftover\n"); asm(" endloop\n"); src += step_row; dst += step_row; } } } JSAMPLE* pixel; // If pixel; RawImage refoc_img; allocing(&refoc_img, nx, ny, 3); float w = F/nu/nv; JSAMPLE* pdst; for(x=0; x<nx; ++x){ for(y=0; y<ny; ++y){ src = &refocus_dst[3*(y+ny*x)]; pdst = GetPixel(&refoc_img, x, y); pdst[0] = (JSAMPLE)CLAMP(src[0]*w); pdst[1] = (JSAMPLE)CLAMP(src[1]*w); pdst[2] = (JSAMPLE)CLAMP(src[2]*w); } } } </pre>	<p>Start SSE2 processing</p> <p>Load destination address</p> <p>Load source address</p> <p>Adding</p> <p>Store value to destination address</p> <p>Processing the leftover with plain assembly code</p> <p>JSAMPLE: unsigned char</p> <p>Normalize the image with the number of additions</p>

Figure 6: SSE2 Implementation of the Derived Algorithm

<pre> void add_data_SIMD_misaligned (unsigned int* pdst, unsigned int* psrc, int src_start, int dst_start, int stride) { // add 'stride' bytes from src[src_start] to dst[dst_start] vec_ushort8 *vdst, *vsrc; int i, loop; vec_ushort8 temp; unsigned int dst = (unsigned int) pdst; unsigned int src = (unsigned int) psrc; if (stride%16){ // stride is odd loop = (int)(stride/16)+1; vec_uchar16 mask = (vec_uchar16)(0x08,0x09,0x0a,0x0b,\ 0x0c,0x0d,0x0e,0x0f,\ 0x10,0x11,0x12,0x13,\ 0x14,0x15,0x16,0x17); if (src_start%2){ // src_start is odd, dst_start is even vsrc = (vec_ushort8*)(src+src_start*8+8); vdst = (vec_ushort8*)(dst+dst_start*8); } else { // src_start is even, dst_start is odd vsrc = (vec_ushort8*)(src+src_start*8); vdst = (vec_ushort8*)(dst+dst_start*8-8); } for (i=0; i<loop; i++){ temp = spu_shuffle(vsrc[i], vsrc[i+1], mask); vdst[i] = spu_add(vdst[i], temp); } } else { loop = (int)(stride/16); vsrc = (vec_ushort8*)(src+src_start*8+16); vdst = (vec_ushort8*)(dst+dst_start*8); for (i=0; i<loop; i++){ vdst[i] = spu_add(vdst[i], vsrc[i]); } } } </pre>	<p>Here we do the main computation</p> <p>pdst: destination address in refocused image</p> <p>psrc: source address in sub-aperture image</p> <p>src_start: the place in psrc to start for each column</p> <p>dst_start: the place in pdst to start for each column</p> <p>stride: number of bytes to be processed per column</p> <p>SPU is vector based. It can process 16 bytes as a vector</p> <p>Care must be taken for data that are not aligned or not a multiple of 16 chars</p> <p>mask: a mask to shuffle the unaligned data. Using the mask, spu_shuffle produces a permuted new vector from two input vectors</p> <p>In this example, a vector contains exactly 2 pixels's color information, therefore, at most one of the source or destination array is not aligned</p> <p>If source is not aligned</p> <p>Adjust the beginning of the source and destination address for shuffling</p> <p>If destination is not aligned</p> <p>Adjust the beginning of the source and destination address for shuffling</p> <p>Shuffle the unaligned data and add up</p> <p>if everything is aligned, it is straightforward to compute the sum</p>
--	--

Figure 8: continued: SPU core function of the CBE Implementation

<pre> void Refocus() int canvas_id, fbo_id, f=0; for(int u=0; u<UV_RES; ++u){ for(int v=0; v<UV_RES; ++v){ g_offsets[u][v] = ((float)u)*beta; g_offsets[v][u] = ((float)v)*beta; } } gClearColor(1.f,1.f,0.f,1); cgGLBindProgram(vertexProgram); cgGLEnableProfile(vertexProfile); cgGLSetStateMatrixParameter(g_CParam_ModelViewProj, CG_GL_MODELVIEW_PROJECTION_MATRIX, CG_GL_MATRIX_IDENTITY); cgGLBindProgram(fragmentProgram); cgGLEnableProfile(fragmentProfile); cgGLSetParameter2f(g_Cfrag_sub_size, RENDERBUFFER_WIDTH\ RENDERBUFFER_HEIGHT); cgGLSetParameter1f(g_Cfrag_total_subs, UV_RES*UV_RES); glViewport(0,0,RENDERBUFFER_WIDTH, \ RENDERBUFFER_HEIGHT); gMatrixMode(GL_MODELVIEW); glLoadIdentity(); glTranslatef(0.f, 0.f, -500.0f); glScalef(RENDERBUFFER_WIDTH/2, RENDERBUFFER_HEIGHT/2, 1.f); canvas_id=0; gBindFramebufferEXT(GL_FRAMEBUFFER_EXT, g_framebuffer); glFramebufferTexture2DEXT(GL_FRAMEBUFFER_EXT, GL_DEPTH_ATTACHMENT_EXT, GL_RENDERBUFFER_EXT, g_depthRenderBuffer); gClear(GL_COLOR_BUFFER_BIT); gBindFramebufferEXT(GL_FRAMEBUFFER_EXT, 0); for(int pass=0; pass<g_lex_numberpass; ++){ fbo_id = canvas_id%2; canvas_id++; canvas_id%=2; gBindFramebufferEXT(GL_FRAMEBUFFER_EXT, \ g_framebuffer); glFramebufferTexture2DEXT(GL_FRAMEBUFFER_EXT, \ </pre>	<pre> GL_COLOR_ATTACHMENT_EXT, \ GL_TEXTURE_RECTANGLE_EXT, \ g_fbo_textid(fbo_id, 0); glFramebufferRenderbufferEXT(GL_FRAMEBUFFER_EXT, GL_DEPTH_ATTACHMENT_EXT, \ GL_RENDERBUFFER_EXT, \ g_depthRenderBuffer); cgGLSetTextureParameter(g_Cfrag_canvas, \ g_fbo_textid(canvas_id)); cgGLEnableTextureParameter(g_CCfrag_canvas); cgGLSetTextureParameter(g_CCfrag_lex, g_fbo_textid(pass)); cgGLEnableTextureParameter(g_CCfrag_lex); cgGLSetParameterArray2f(g_CCfrag_offsets, 0, SHADER_LOOP, &g_offsets[pass*SHADER_LOOP*2]); CheckGLError(); drawImage(); gBindFramebufferEXT(GL_FRAMEBUFFER_EXT, 0); } cgGLDisableTextureParameter(g_CCfrag_canvas); cgGLDisableTextureParameter(g_CCfrag_lex); cgGLDisableProfile(vertexProfile); cgGLDisableProfile(fragmentProfile); } //Fragment Shader void main(in float2 canvas_coord: TEXCOORD, uniform samplerRECT canvas_lex: TEXTUNIT0, uniform samplerRECT lex: TEXTUNIT1, uniform float2 offsets[4], uniform float2 sub_size, uniform float total_subs, out float3 OUT : COLOR) { float2 temp=float2(0,0); float u,v; float2 coord; float n=0; for(u=0; u<2; ++u){ for(v=0; v<2; ++v){ coord = float2(v*sub_size, u \ *sub_size); u*sub_size>canvas_coord+float2(offsets[n],y, offsets[n],x); temp += texRECT(lex, coord); n++; } } OUT = temp/total_subs+RECT(canvas_lex, canvas_coord); } </pre>
<p>Compute the offset for each sub-aperture image</p> <p>Preparation for multi-pass rendering. According to GPU capacity, a maximum number of sub-aperture images are tiled into one texture. In each pass, pixels are extracted from the texture and added with the canvas. The canvas is a texture associated with frame buffer, which stores aggregating results during the multi-pass. We use 2 frame buffer textures to achieve this.</p> <p>Bind CG programs and pass parameters to it.</p> <p>Set Viewport to be the same size of the refocused image</p> <p>We essentially render a square with the refocused texture on it.</p> <p>To do that, we need some transformation so that the square fits well within the viewport</p> <p>Initialize the frame buffer and clear it with 0 (black color)</p> <p>Start the multipass</p> <p>The result of the current pass is stored in the frame buffer texture. In the next pass, this texture serves as a canvas, then it is added with highfield data stored in more sub-aperture images. The result is stored in the other frame buffer texture, which becomes the</p>	<p>canvas in the next pass.</p> <p>Bind CG program (the fragment shader)</p> <p>Pass the parameters</p> <p>Draw the square with the CG fragment shader</p> <p>A CG fragment shader implements the refocusing algorithm with 2*2 sub-aperture images</p> <p>canvas_coord: rasterized coordination of canvas</p> <p>canvas_lex: The texture for recorded refocusing process by far lex. The texture for the new set of sub-aperture images</p> <p>offsets: The offsets of the 4 sub-aperture images</p> <p>sub_size: The size of each sub-aperture image</p> <p>total_subs: Number of all the sub-aperture images, used for normalization</p> <p>Compute the coordinates for each pixel and add them up</p> <p>Normalize the color data and add it to the canvas color</p>

Figure 9: CG fragment shader for the Derived Algorithm