

Parallel Implementation of a Visual Attention Model

Anis RAHMAN

June 28, 2010

Table of contents

- 1 Overview
 - Vision systems
 - Case study: Visual saliency model
 - Characteristics of such applications
 - Why use GPU as co-processor?
 - Motivation
- 2 GPU
 - The programming model(1)
 - The programming model(2)
- 3 Optimizations
 - How to program?
 - Other options
- 4 Issues
 - Issues
- 5 Conclusion

Overview

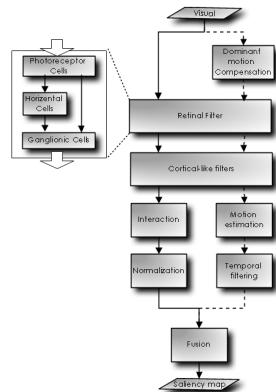
- GPUs and vision systems
- Case study: Visual attention model

Vision systems

- Sub-field of Computer vision
- To build artificial systems
- Involves large sets of data
- Same set of operations on it

Case study: Visual saliency model

- To find the spotlight of focus
- Based on human visual system
- Bottom-up model
- Implements both pathways
- Some applications:
 - To automate cinematography, surveillance, and video reframing
 - To simulate mediated reality
 - To find ROI maps



Saliency map



(a) Original image



(b) M_s



(c) M_d



(d) Saliency map

Characteristics of such applications

- Storage and memory usage
- Performance
 - Requires real time capability
 - Involves complex computations
- Data parallelism
 - Single operation on huge data
 - No or less dependency
- GPU friendly operations
- Previous attempts made required:
 - Learning graphics specific APIs
 - Re-structuring of algorithms according to the graphics pipeline

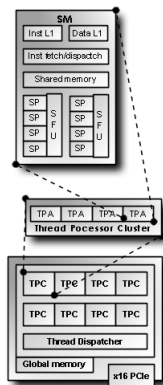
Why use GPU as co-processor?

- Stream processing architecture
 - All processors work in groups
 - Can communicate through shared memory
- Performance is unmatched
- Very high memory bandwidth
- Accessible
- Easier to program and manage
- Already applied in diverse fields
 - Biological engineering
 - Oil and gas exploration
 - Financial analysis

- To port data-parallel vision algorithm
- To demonstrate the speedup
- To confirm effects of low precision
- To apply different optimizations
- To experience the difficulties

The programming model

- CUDA(Compute Unified Device Architecture)
 - Independent of traditional graphics APIs
 - Based on:
 - Arithmetic intensity = Arithmetic / Bandwidth
 - Uses familiar C
- Allows access to on-chip shared memory
- Provides texture lookups
- Supported by GPU-specific libraries
 - CUFFT, CUDPP, CuBLAS, OpenVIDIA



The programming model

- Code is composed of:
 - Host code
 - Kernel code
- Two-level thread hierarchy
 - 2D grid of thread blocks
 - Each thread block is 3D grid of threads
- Making the code scalable

How to program?

Restructuring the algorithm

- Identifying the data-parallel portions
- Avoiding code divergence

Effective use of memory model

- Constant cache for persistent data values
- Texture cache for frequently accessed data values
- Shared memory to reduce multiple accesses to global memory

Other options

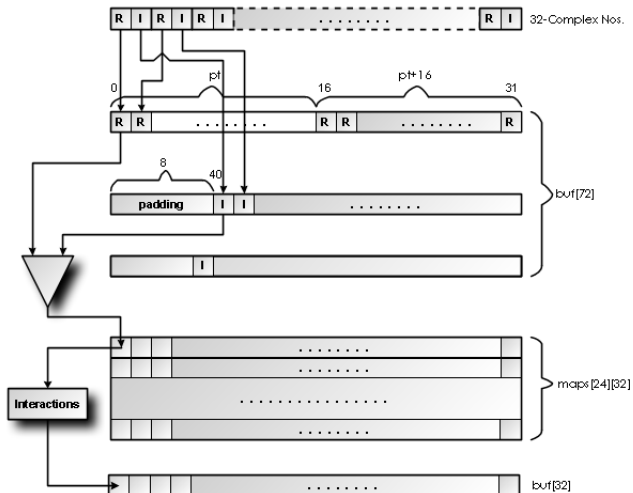
Using NVCC compiler options

- `-use_fast_math`: to use CUDA fast math functions
- `-arch sm_13`: to enable double precision (if supported)

Using GPU-specific libraries

- GPU-specific functions
- Black-box algorithm

Example: Interactions kernel



```

1  __global__ void ShortInteractionKernel ( Complex* in, unsigned int width
2      , unsigned int height, float* out) {
3      __shared__ float maps[ NO_OF_ORIENTS * NO_OF_BANDS ][32];
4      __shared__ float buf [72];
5
6      unsigned int x1 = blockIdx.x*blockDim.x + threadIdx.x/2;
7      unsigned int x2 = blockIdx.x*blockDim.x + threadIdx.x;
8      unsigned int y  = blockIdx.y*blockDim.y + threadIdx.y;
9
10     if ( x1 >= width || x2 >= width || y>= height) return;
11
12     unsigned int mod  = threadIdx.x%2;
13     unsigned int pt   = threadIdx.x/2 + 40* mod;
14     unsigned int size = width*height;
15
16     for ( unsigned int j=0 ; j< NO_OF_ORIENTS ; ++j) {
17         for ( unsigned int i=0 ; i< NO_OF_BANDS ; ++i) {
18             /* *****
19              * 32 threads process 16 complex numbers in parallel
20              * every thread stores them with real and imaginary interlaced
21              * 32 threads produce 32 real products in parallel
22              ***** */
23             buf[pt] = // first 16 complex numbers
24                 in[(j* NO_OF_BANDS+i)*size+(y*width+x1)][mod]/(float)(size);
25             buf[pt+16] = // next 16 complex numbers
26                 in[(j*NO_OF_BANDS+i)*size+(y*width+x1+16)][mod]/(float)(size);
27             __syncthreads();
28
29             maps[j* NO_OF_BANDS + i][ threadIdx.x] = abs(
30                 buf[threadIdx.x ]*buf[threadIdx.x ] +
31                 buf[threadIdx.x + 40]* buf[threadIdx.x + 40]);
32             __syncthreads ();
33         }
34     }
35     // prefetched data in shared memory is used by interactions
36 }

```

Speedup

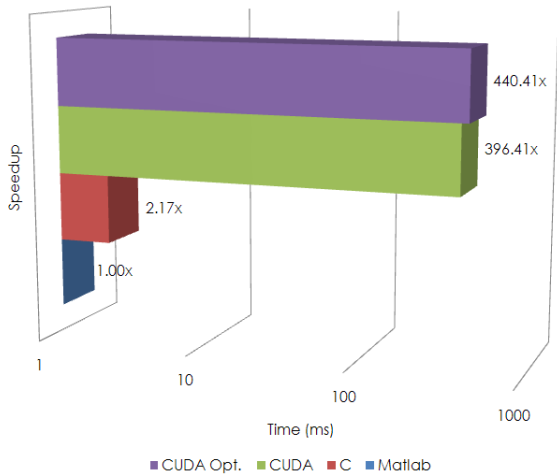


Figure: Speedup for the visual saliency model

Table: Computational cost of each step in static pathway

Kernel	Geforce GTX 260 (ms)	Geforce 8800 GTS (ms)
Mask	0.08	0.57
FFT	0.59	1.36
Shift	0.09	0.21
$24 \times$ Gabor	1.47	6.31
$24 \times$ Inverse shift	1.13	2.66
$24 \times$ IFFT	10.76	32.54
$24 \times$ Interaction	3.13	7.06
$24 \times$ Normalize	3.33	25.27
$24 \times$ Normalize Itti	3.34	25.37
$24 \times$ Normalize Fusion	2.89	21.90
Total	26.81	123.24

Table: Computational cost of each step in dynamic pathway

Kernel	Geforce GTX 285 (ms)
Retinal Filtering	21.5
Modulation	2.6
Demodulation	3.1
Interpolation	0.3
Projection	0.2
Ver. Guassian recursive	33.2
Hor. Guassian recursive	21.7
Gradients	6.2
MCPI	39.1
Median filtering	0.2
Total	128.1

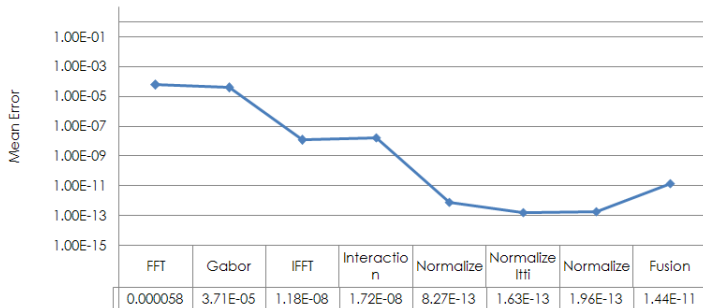


Figure: The effect of lower precision support on the result

Universal image index = $Q = 99.66\%$

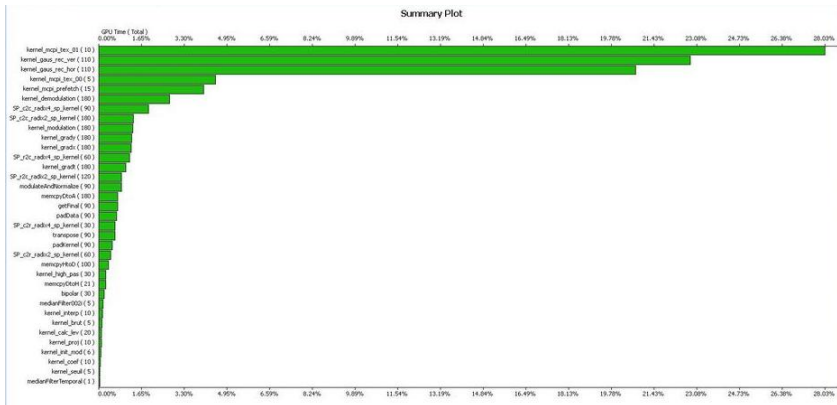


Figure: Profiling graph for the dynamic path

Multi-GPU solution

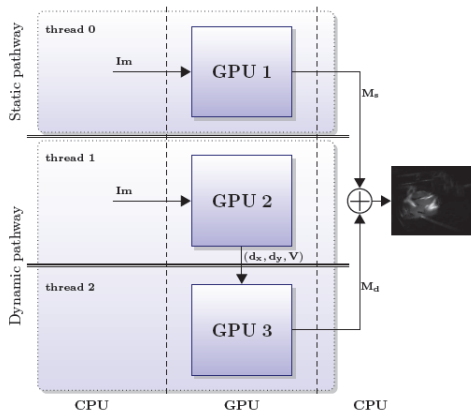


Figure: Block diagram of multi-GPU model

Issues

Memory size

- Only one image loaded to GPU, due to lack of GPU memory
- Causes frequent context switching
- CPU-GPU memory transfer i.e. one image per pass

Precision

- Doesn't fully conform to IEEE-754 standard
- Yet no full support for double precision
- Can lead in unusable results, due to lower accuracy
- Lower precision can be resolved using mixed precision

Conclusion

- Real time capability (~ 28 fps)
- Created opportunity to extend the model
- Exploited GPUs power without extensive re-structuring
- Without the need for high precision