# Image Processing Optimization C# on GPU with Hybridizer™

regis.portalez@altimesh.com



# Median Filter Denoising



Noisy image (lena 1960x1960)



Denoised image window = 3



Median Filter Denoising



For each pixel, we read  $(2 * window + 1)^2$  pixels of input

### Optimization Steps An Overview

<ol> <li>Enable C# parallelization (remove loop side effects)</li> <li>Use Parallel.For</li> </ol>	Necessary	x5
3. Run on GPU (Hybridizer)		
- 3.1 Decorate methods		_
– 3.2 Allocate memory	Low cost	x78
– 3.3 Feed our 50k threads		
4. Implement Advanced Optimizations		
- 4.1 Shared memory	Bonus	x92
- 4.2 Texture memory		
5. More Optimizations	Expertise	x?

Median Filter is not easy. On easier code, steps 3 and 4 would be sufficient

#### AForge code

```
ushort* src, dst;
for (int y = startY; y < stopY; y++)</pre>
{
    for (int x = startX; x < stopX; x++, src++, dst++)</pre>
    {
         int c = 0;
         for (i = -radius; i <= radius; i++)</pre>
         {
             for (j = -radius; j <= radius; j++)</pre>
             {
                 g[c++] = src[i * srcStride + j];
             }
         }
         Array.Sort(g, 0, c);
         *dst = g[c >> 1];
    }
    src += srcOffset;
    dst += dstOffset;
}
```



#### AForge code



Old-school optimizations Inner loops have sideeffects Requires unsafe



# Enable Parallelization Remove loop side-effects

```
var buffer1 = new ushort[windowCount * windowCount];
for (int j = window; j < height - window; ++j)</pre>
{
    for (int i = window; i < width - window; ++i)</pre>
    {
        for (int k = -window; k <= window; ++k)</pre>
        {
             for (int p = -window; p <= window; ++p)</pre>
             {
                 int bufferIndex = (k + window) * windowCount + p + window;
                 int pixelIndex = (j + k) * width + (i + p);
                 buffer1[bufferIndex] = input[pixelIndex];
             }
        }
        Array.Sort(buffer1, 0, windowCount * windowCount);
        output[j * width + i] = buffer1[(windowCount * windowCount) / 2];
    }
}
```





No performance penalty – Jitter is quite smart now! Much more readable code Inner loops are independant of the outer loops: possible to introduce parallelization

#### 2. Use Parallel.For

```
Parallel.For(window, height - window, j =>
{
    var buffer1 = new ushort[windowCount * windowCount];
    for (int i = window; i < width - window; ++i)</pre>
    {
        for (int k = -window; k <= window; ++k)</pre>
        {
            for (int p = -window; p <= window; ++p)</pre>
            {
                 int bufferIndex = (k + window) * windowCount + p + window;
                 int pixelIndex = (j + k) * width + (i + p);
                 buffer1[bufferIndex] = input[pixelIndex];
            }
        }
        Array.Sort(buffer1, 0, windowCount * windowCount);
        output[j * width + i] = buffer1[(windowCount * windowCount) / 2];
    }
});
```





One line change yields a x5,5 speed-up



3. Run On GPU CUDA & GPU: A Few Words



- Multiprocessors (SM) are
- similar to CPU Cores
- CUDA cores are similar to CPU SIMD lanes



# 3. Run On GPU CUDA threading model



- Threads are grouped in blocks
- Blocks are grouped in a grid
- Grids and blocks have configurable shape (1, 2 or 3D)
- I block run on a single SM



### 3. Run on GPU Hybridizer™: A Few Words

- Hybridizer™ is a compiler targeting CUDA-enabled GPUS from DotNet.
- Attribute-based (no runtime cost)

• ...

- Integrated with debugger and profiler
- Support of Generics and Virtual functions

- Trial version downloadable from <u>Visual Studio Marketplace</u>
- Professional edition available in beta (<u>Altimesh website</u>)
- Full version already deployed in Investment Banks (upon request)

### 3.1 Run On GPU Decorate Methods

```
One and only
               [EntryPoint]
modification
               public static void ParallelCsharp(byte[] output, byte[] input, int width, int height)
               {
                    Parallel.For(window, height - window, j =>
                    {
                        var buffer1 = new byte[windowCount * windowCount];
                        for (int i = window; i < width - window; ++i)</pre>
                        {
                            for (int k = -window; k <= window; ++k)</pre>
                            {
                                for (int p = -window; p <= window; ++p)</pre>
                                    int bufferIndex = (k + window) * windowCount + p + window;
                                    int pixelIndex = (j + k) * width + (i + p);
                                    buffer1[bufferIndex] = input[pixelIndex];
                                }
                            }
                            Array.Sort(buffer1, windowCount * windowCount);
                            output[j * width + i] = buffer1[(windowCount * windowCount) / 2];
                        }
                    });
               }
```





#### Relative Performance

Quite disappointing isn't it? WHY??



### 3.2 Allocate Memory Heap Allocation On GPU





#### 3.2 Allocate Memory Move To Stack

```
[EntryPoint]
public static void ParallelCsharp(ushort[] output, ushort[] input, int width, int height)
{
    Parallel.For(window, height - window, j =>
    {
        var buffer1 = new StackArray<ushort>(windowCount * windowCount);
        for (int i = window; i < width - window; ++i)
        {
            …
        }
        });
}</pre>
```



#### Relative Performance





#### 3.3 Feed the Beast

#### CPU

- Cores
  - Consumer : 8
  - Server : 22
- SIMD Lanes
  - AVX2 : 4 8
  - AVX512 : 8 16
- Hyperthreading
   x2
- Parallelism



32 up to 704

#### GPU

- SMs
  - GeForce : 28
  - Tesla : 80
- Cores per SM - GeForce : 128
  - Tesla : 64
- Context (to hide latency)
   32
- Parallelism



3,584 up to 164,000







#### 3.3 Feed the Beast Use A 2D Grid



#### Relative Performance



Run time (seconds):

- AForge : 4,16
- Parallel C# : 0,76
- Hybridizer Stack 2D : 0,053

Can we do better?



Run time (seconds):

- AForge : 4,16
- Parallel C# : 0,76
- Hybridizer Stack 2D : 0,053

Can we do better?



Seems we are reading too much data!

4.1 Implement Advanced Optimizations Leverage On-Chip Cache (Shared memory)



Common read zone



4.1 Implement Advanced Optimizations Leverage On-Chip Cache (Shared memory)

window

Block Should be cached



Common read zone



4.1 Implement Advanced Optimizations Leverage On-Chip Cache (Shared memory)



- On chip (Multiprocessor)
- Accessible by entire block
- 48KB per block
- See it as CPU L1-cache with explicit control



#### 4.1 Implement Advanced Optimizations Leverage On-Chip Cache (Shared memory) [EntryPoint] public static void Parallel2DShared(ushort[] output, ushort[] input, int width, int height) { int cacheWidth = blockDim.x + 2 \* window; Cache « allocation » ushort[] cache = new SharedMemoryAllocator<ushort>().allocate(cacheWidth\* cacheWidth\*); for (int bid\_j = blockIdx.y; bid\_j < (height ) / blockDim.y; bid\_j += gridDim.y)</pre> { for (int bid\_i = blockIdx.x; bid\_i < (width) / blockDim.x; bid\_i += gridDim.x)</pre> { int bli = bid i \* blockDim.x; int blj = bid\_j \* blockDim.y; int i = threadIdx.x + bid i \* blockDim.x; int j = threadIdx.y + bid\_j \* blockDim.y; // ... some code to fetch cache - put data in shared memory Synchronize threads CUDAIntrinsics.\_\_syncthreads(); < in block var buffer1 = new StackArray<ushort>(windowCount \* windowCount); var buffer2 = new StackArray<ushort>(windowCount \* windowCount); for (q = -window; q <= window; ++q)</pre> { for (p = -window; p <= window; ++p)</pre> { int bufferIndex = (q + window) \* windowCount + p + window; int cacheIndex = (threadIdx.y + window + q) \* cacheWidth + threadIdx.x + window + p; Read from cache buffer1[bufferIndex] = cache[cacheIndex]; } } MergeSort(buffer1, buffer2, windowCount \* windowCount); output[j \* width + i] = buffer1[(windowCount \* windowCount) / 2]; } } }





#### Relative Performance

We have a x87 speed-up over initial single-threaded code.

Code still works in .Net

Can we do better?



ถ

# 4.2 Implement Advanced Optimizations Leverage Texture Cache



- Different memory cache
- Optimized for 2D spatial locality





# 4.2 Implement Advanced Optimizations Leverage Texture Cache



- Different memory cache
- Optimized for 2D spatial



Bind input image to texture

### 4.2 Implement Advanced Optimizations Leverage Texture Memory

IntPtr src = runner.Marshaller.MarshalManagedToNative(pixels);

cudaChannelFormatDesc channelDescTex = TextureHelpers.cudaCreateChannelDesc(8, 0, 0, 0, cudaChannelFormatKind.cudaChannelFormatKindUnsigned); cudaArray\_t cuArrayTex = TextureHelpers.CreateCudaArray(channelDescTex, src, width, height); cudaResourceDesc resDescTex = TextureHelpers.CreateCudaResourceDesc(cuArrayTex);

//create Texture descriptor
cudaTextureDesc texDesc = TextureHelpers.CreateCudaTextureDesc();

//create Texture object
cudaTextureObject\_t texObj;
cuda.CreateTextureObject(out texObj, ref resDescTex, ref texDesc);

- CUDA API is fully available through a wrapper (P/Invoke)
- Texture and Surface API types are exposed and mapped (IntrinsicTypes)
- Resulting C# code for textures usage very similar to CUDA/C tutorials





#### **Relative Performance**

# We accelerated AForge with a x92 speed-up.

Can we do better?







# 5. Implement Advanced Optimizations What's next?



Put everything in register file

# GPU SM have 32k registers for a Block - up to 255 by threads



5.1 Implement Advanced Optimizations Rolling Buffer Of Registers

	(i,j)	
	(i,j+1)	

Load data in registers and process pixel (i,j)



5.1 Implement Advanced Optimizations Rolling Buffer Of Registers



Load next line and roll buffer for pixel(i, j+1)



#### 5.2 Implement Advanced Optimizations Loop Unrolling



If window is a compile-time constant, backend-compiler is able to completely unroll loop

(actually required for compiler to map arrays on registers)

#### 5.3 Implement Advanced Optimizations Smart Sorting

- Sorting networks are optimal for known-size arrays.
- They are not capable of sorting arbitrary long arrays.
- Possible to implement in C++ meta-programming.
- Enabled with hand-written CUDA, called from C# using « IntrinsicType »

```
template <typename scalar, int window>
                                                              struct medianfilter
                                                                static constexpr int size = (window * 2 + 1) * (window * 2 + 1);
                                                                scalar buffer[size];
                                                                scalar work[size];
[IntrinsicInclude("intrinsics.cuh")]
                                                               _forceinline_ __device_ __host__ void set_Item(int i, scalar val) {
[IntrinsicType("medianfilter<unsigned short, 3>")]
                                                              buffer[i] = val; }
struct medianfilter ushort 3
                                                              __forceinline___device___host__ scalar apply() {
{
    public ushort apply() { ... }
                                                                #pragma unroll
    public void rollbuffer() { ... }
                                                                for (int k = 0; k < size; ++k) {</pre>
    public ushort get Item(int i) { ... }
                                                                  work[k] = buffer[k];
    public void set_Item(int i, ushort val) { ... }
                                                                }
}
                                                                hybridizer::StaticSort<size> sort;
                                                                sort(work);
                                                                return work[size / 2];
                                                              }
```



#### **Relative Performance**



### Can we do better?



#### 6. Write Plain CUDA

• Writing the entire application in CUDA/C leads to



#### Can we do better?



Maybe...



We barely read the image once

- 7.68 MB read
- 7.68 MB write
- 0.49 MB overhead



### Maybe...





#### 7. Take Away



CPU: 1.8y ACC: 1.9y

BANDWIDTH doubles every

CPU: 4.3y ACC: 2.8y

Caching computations is not necessary anymore

Caching memory operations is mandatory!

Always use the fastest memory available, the fastest of them all being registries



# Memory interaction is the elephant in the room

# Thank you



#### **Relative Performance**



#### http://www.altimesh.com

Windows 10 x64

