# Parallel Image Processing with CUDA
## – A case study with the Canny Edge Detection Filter –

Daniel Weingaertner

Informatics Department
Federal University of Paraná - Brazil

Hochschule Regensburg
02.05.2011

# Summary

# Paraná – Brazil
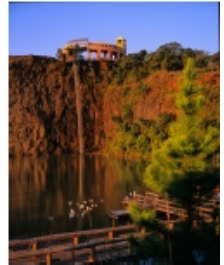
# Brazil – Europe

# Paraná

# Curitiba

# Federal University of Paraná



Informatic Department

# Informatics Department

Undergraduate: Bachelor in Computer Science

- 8 semesters course
- 80 incoming students per year

Bachelor in Biomedical Informatics

- 8 semesters course
- 30 incoming students per year

Graduate: Master and PhD in Computer Science

- Algorithms, Image Processing, Computer Vision, Artificial Intelligence
- Databases, Scientific Computing and Open Source Software, Computer-Human Interface
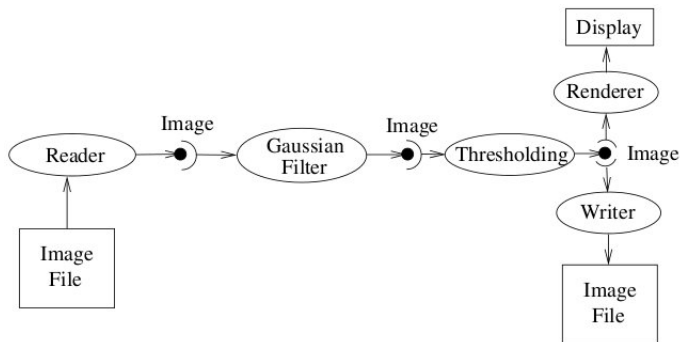- Computer Networks, Embedded Systems

# Summary

# Insight Toolkit (ITK)

- Created in 1999, Open Source, Multi platform, Object Oriented (Templates), Good documentation and support



**Figure:** Image Processing Workflow in ITK

# ITK - Sample code

```
1  #include "itkImage.h"
2  #include "itkImageFileReader.h"
3  #include "itkImageFileWriter.h"
4  #include "itkCannyEdgeDetectionImageFilter.h"
5
6  typedef itk::Image<float,2>                     ImageType;
7  typedef itk::ImageFileReader< ImageType > ReaderType;
8  typedef itk::ImageFileWriter< ImageType > WriterType;
9  typedef itk::CannyEdgeDetectionImageFilter< ImageType, ImageType > CannyFilter;
10
11 int main (int argc, char** argv){
12
13   ReaderType::Pointer reader = ReaderType::New();
14   reader->SetFileName( argv[1] );
15   reader->Update();
16
17   CannyFilter::Pointer canny = CannyFilter::New();
18   canny->SetInput(reader->GetOutput());
19   canny->SetVariance( atof( argv[3] ) );
20   canny->SetUpperThreshold( atoi( argv[4] ) );
21   canny->SetLowerThreshold( atoi( argv[5] ) );
22   canny->Update();
23
24   WriterType::Pointer writer = WriterType::New();
25   writer->SetFileName( argv[2] );
26   writer->SetInput( canny->GetOutput() );
27   writer->Update();
28
29   return EXIT_SUCCESS;
30 }
```
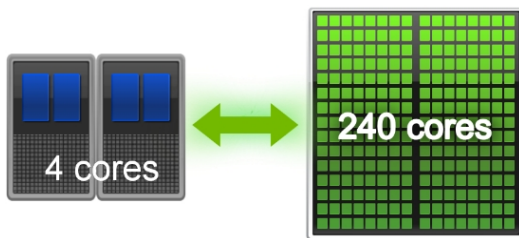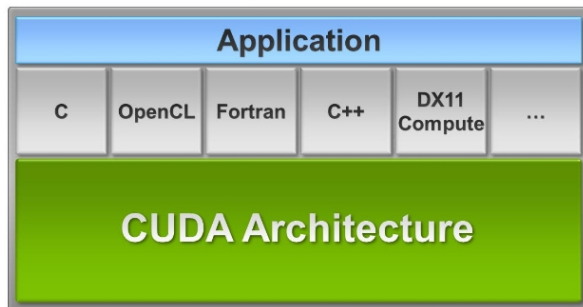
# Summary

# What is GPGPU Computing?

- The use of the GPU for general purpose computation
- CPU and GPU can be used concurrently
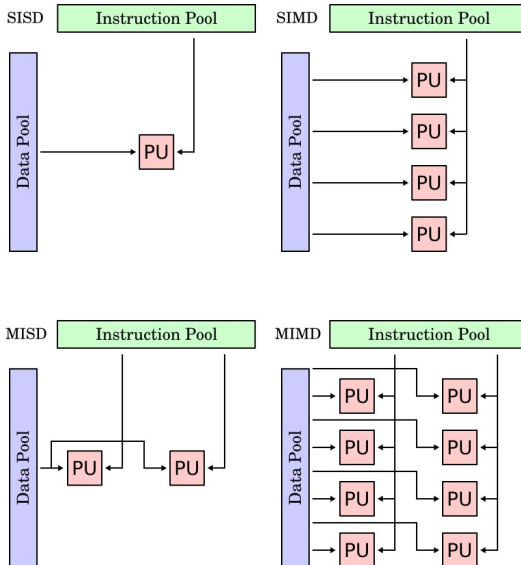- To the end user, its simply a way to run applications faster.

# What is CUDA?

- CUDA = *Compute Unified Device Architecture*.
- *General-Purpose Parallel Computing Architecture*.
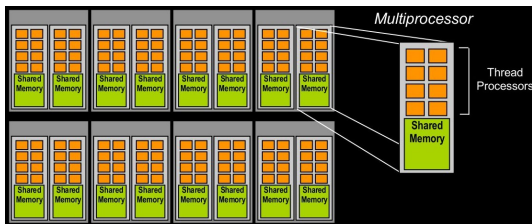- Provides libraries, C language extension and hardware driver.
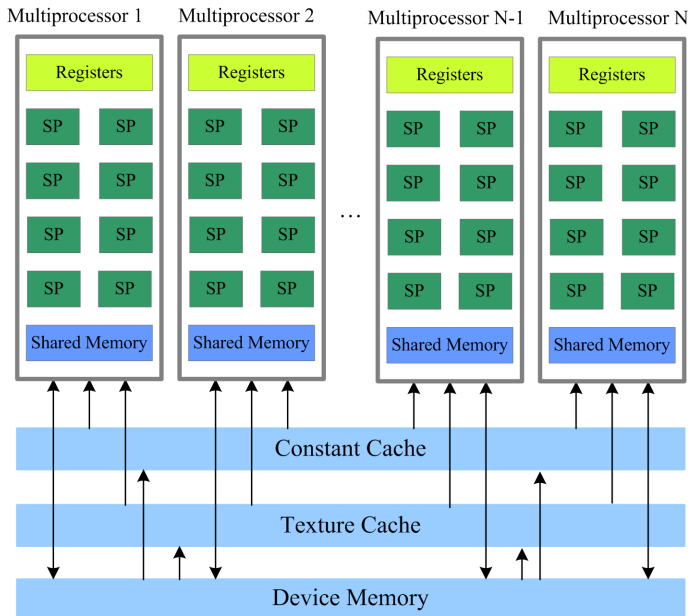
# Parallel Processing Models

# Single-Instruction Multiple-Thread Unit

- Creates, handles, schedules and executes groups of 32 threads (*warp*).
- All threads in a *warp* start at the same point.
- But they are "free" to jump to different code positions independently.

# CUDA Architecture Overview

# Optimization Strategies for CUDA

**Main optimization strategies for CUDA involve:**

- Optimized/careful memory access
- Maximization of processor utilization
- Maximization of non-serialized instructions

## CUDA - Sample Code

```c
1  #include <stdio.h>
2  #include <assert.h>
3  #include <cuda.h>
4  void incrementArrayOnHost(float *a, int N)
5  {
6    int i;
7    for (i=0; i < N; i++) a[i] = a[i]+1.f;
8  }
9  __global__ void incrementArrayOnDevice(float *a, int N)
10 {
11   int idx = blockIdx.x*blockDim.x + threadIdx.x;
12   if (idx<N) a[idx] = a[idx]+1.f;
13 }
14 int main(void)
15 {
16   float *a_h, *b_h;            // pointers to host memory
17   float *a_d;                  // pointer to device memory
18   int i, N = 10000;
19   size_t size = N*sizeof(float);
20   a_h = (float *)malloc(size);
21   b_h = (float *)malloc(size);
22   cudaMalloc((void **) &a_d, size);
23   for (i=0; i<N; i++) a_h[i] = (float)i;
24   cudaMemcpy(a_d, a_h, sizeof(float)*N, cudaMemcpyHostToDevice);
25   incrementArrayOnHost(a_h, N);
26   int blockSize = 256;
27   int nBlocks = N/blockSize + (N%blockSize == 0?0:1);
28   incrementArrayOnDevice <<< nBlocks, blockSize >>> (a_d, N);
29   cudaMemcpy(b_h, a_d, sizeof(float)*N, cudaMemcpyDeviceToHost);
30   free(a_h); free(b_h); cudaFree(a_d);
31 }
```

# Summary

# Integrating CUDA Filters into ITK Workflow

**ITK community suggests:**

- Re-implement filters where parallelizing provides significant speedup
- Consider the entire workflow: copying to/from the GPU is very time consuming

**Careful!**

"Premature optimization is the root of all evil!" (Donald Knuth)

# Integrating CUDA Filters into ITK Workflow

**ITK community suggests:**

- Re-implement filters where parallelizing provides significant speedup
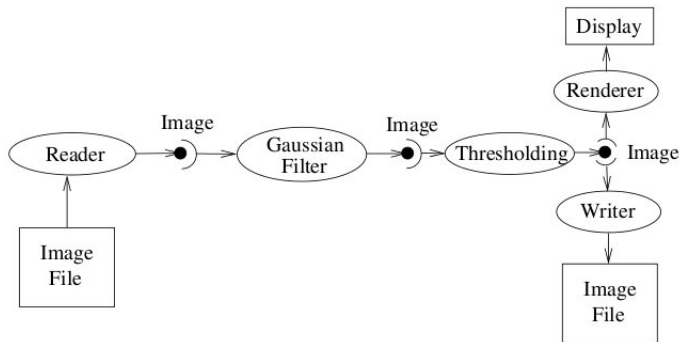- Consider the entire workflow: copying to/from the GPU is very time consuming

**Careful!**

"Premature optimization is the root of all evil!" (Donald Knuth)

# CUDA Insight Toolkit (CITK)

**Changes to ITK**

- Slight architecture change: CudaImportImageContainer
- Backwards compatible
- Data transfer between HOST and DEVICE only "on demand"
- Allows for filter chaining inside the DEVICE

# Summary

# CudaCanny

- itkCudaCannyEdgeDetectionImageFilter

**Algorithm 1** Canny Edge Detection Filter

Gaussian Smoothing
Gradient Computation
Non-Maximum Supression
Histeresis

# Gradient Computation with Sobel Filter

- itkCudaSobelEdgeDetectionImageFilter

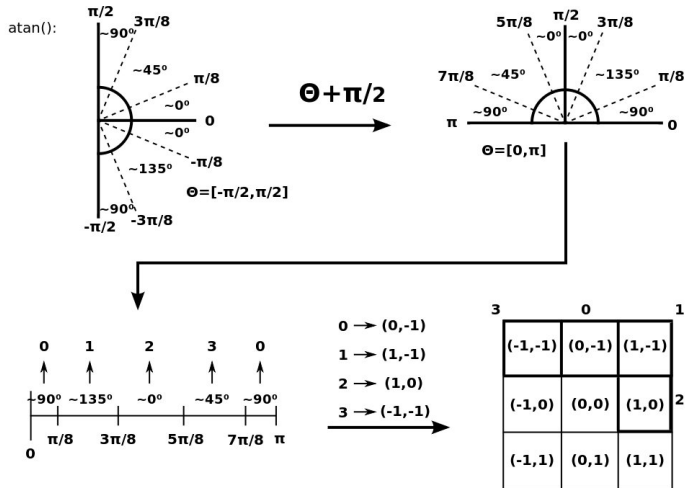| -1 | 0 | 1 |
|---|---|---|
| -2 | 0 | 2 |
| -1 | 0 | 1 |

(a) Sobel X

| 1 | 2 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -2 | -1 |

(b) Sobel Y

$$L_v = \sqrt{L_x^2 + L_y^2} \tag{1}$$

$$\theta = \arctan\left(\frac{L_y}{L_x}\right) \tag{2}$$

# Optimization for Edge Direction Computation

# Code Extract from CudaSobel

```
if ((pos.x) && ((size.x-1)-pos.x) && (pos.y) && ((size.y-1)-pos.y)){

  diagonal.x = tex1Dfetch(texRef,(pixIdx-size.x-1));
  diagonal.y = tex1Dfetch(texRef,(pixIdx-size.x+1));
  diagonal.z = tex1Dfetch(texRef,(pixIdx+size.x-1));
  diagonal.w = tex1Dfetch(texRef,(pixIdx+size.x+1));
  cross.x = tex1Dfetch(texRef,(pixIdx-size.x));
  cross.y = tex1Dfetch(texRef,(pixIdx+size.x));
  cross.z = tex1Dfetch(texRef,(pixIdx-1));
  cross.w = tex1Dfetch(texRef,(pixIdx+1));

  /// SobelX
  g_i.x -= (diagonal.x+cross.z+cross.z+diagonal.z);
  g_i.x += (diagonal.y+cross.w+cross.w+diagonal.w);

  /// SobelY
  g_i.y -= (diagonal.z+cross.y+cross.y+diagonal.w);
  g_i.y += (diagonal.x+cross.x+cross.x+diagonal.y);
}

Magnitude[pixIdx] = sqrtf((g_i.x*g_i.x) + (g_i.y*g_i.y));

theta = (g_i.x != 0)*(int)(atanf(__fdividef(g_i.y,g_i.x))*__fdividef(180,M_PI)) + 90;
if (theta > 157) theta -= 158;
theta = ceilf(__fdividef(theta-22,45));
Direction[pixIdx] = make_short2(1-(theta == 0)-((theta == 1)<<1),(theta == 2)-1);

}
```
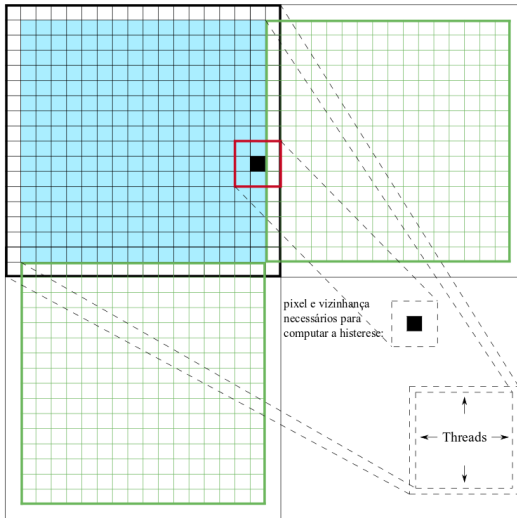
# Histeresis Operation



pixel e vizinhança necessários para computar a histerese

← Threads →

# Histeresis Algorithm

---
**Algorithm 2** Histeresis on CPU

  Transfers the Gradient/NMS images to the GPU
  **repeat**
    Run the histeresis kernel on GPU
  **until** no pixel changes status
  Return edge image

---

# Histeresis Algorithm

**Algorithm 3** Histeresis on GPU

Load an image region with size 18x18 into shared memory
modified $\leftarrow$ false
**repeat**
  modified_region $\leftarrow$ false
  Synchronize *threads* of same multiprocessor
  **if** Pixel changes status **then**
    modified $\leftarrow$ true
    modified_region $\leftarrow$ true
  **end if**
  Synchronize *threads* of same multiprocessor
**until** modified_region $=$ false
**if** modified $=$ true **then**
  Update modified status on HOST
**end if**

# Summary

# Metodology

**Hardware:**

- Server:
  - CPU: 4x AMD Opteron(tm) Processor 6136 2,4GHz with 8 cores, each with 512 KB cache and 126GB RAM
  - GPU1: NVidia Tesla C2050 with 448 1,15GHz cores and 3GB RAM.
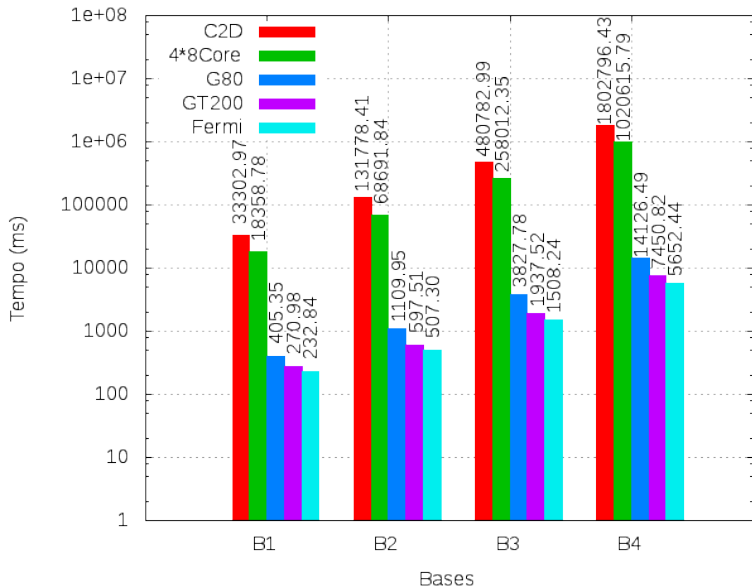  - GPU2: NVidia Tesla C1060 com 240 1,3GHz cores and 4GB RAM.
- Desktop:
  - CPU: Intel®Core(TM)2 Duo E7400 2,80GHz with 3072 KB cache and 2GB RAM
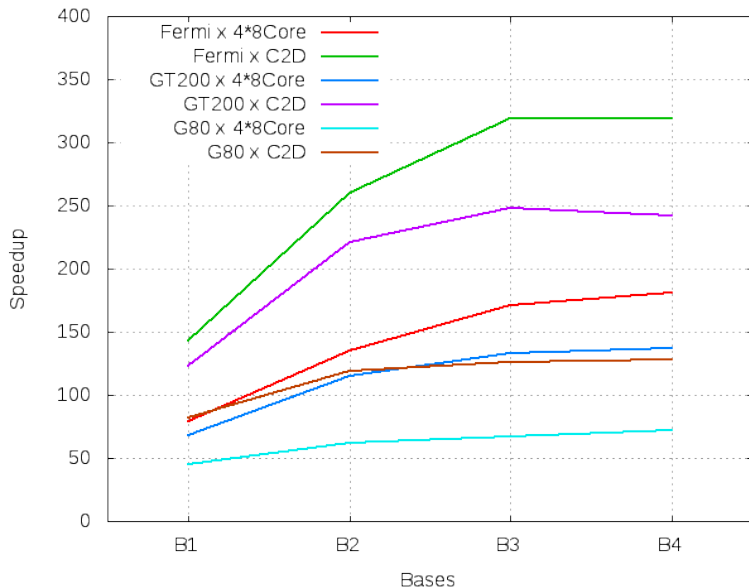  - GPU: NVidia GeForce 8800 GT with 112 1,5GHz cores and 512MB RAM.

# Metodology

- Images from the Berkeley Segmentation Dataset

| Base | Image resolution | Num. of Images |
|------|------------------|----------------|
| B1 | 321×481 e 481×321 | 100 |
| B2 | 642×962 e 962×642 | 100 |
| B3 | 1284×1924 e 1924×1284 | 100 |
| B4 | 2568×3848 e 3848×2568 | 100 |

# Performance Tests

# Performance Tests

# Performance Tests

# Performance Tests

# Summary

# Conclusion

**Parallel Programming**

- Parallel programming is definitely the way to go.
- Implement efficient parallel code is demanding.
- Programmer should know more details about the hardware, especially memory architecture.

**Canny Filter with CUDA**

- We had a great speedup on the edge detection filter
- Also noticed that the existing implementation is not efficient
- There is still a LOT of work if we want to parallelize ITK.

# Conclusion

## Parallel Programming

- Parallel programming is definitely the way to go.
- Implement efficient parallel code is demanding.
- Programmer should know more details about the hardware, especially memory architecture.

## Canny Filter with CUDA

- We had a great speedup on the edge detection filter
- Also noticed that the existing implementation is not efficient
- There is still a LOT of work if we want to parallelize ITK.

# Contact

**Thank You!**

Daniel Weingaertner
danielw@inf.ufpr.br