Parallelization of visual perception model

Anis RAHMAN

GIPSA-Lab, Grenoble Institute of Technology, France email:anis.rahman@gipsa-lab.grenoble-inp.fr

PhD supervised by: Dominique Houzet & Denis Pellerin

Team: AGPIG (Architecture Géométrie, Perception, Images, Gestes)

June 28, 2010



Introduction Parallel implementation Results Conclusion

Outline



Introduction

- 2 Parallel implementation
 - Porting onto GPU
 - Multi-GPU platform
- 3 Results
 - Performance gains
 - Effects of lower precision

Onclusion

Visual perception

"Visual perception is the ability to interpret information and surroundings from visible light reaching the eye."

Modeling approach

- Forward: by study of natural system
- **Reverse:** by building artificial system

Visual perception

Importance to model vision

- to understand the properties of natural vision
- to create an efficient and robust vision system
- to explore its potential for numerous applications

Importantly, it will be a step towards *building brains*.

Towards a robust vision system

Modeling issues

- involve a lots of parameters to test i.e. large number of neurons
- are complex and compute-intensive
- are difficult to study and to improve

Towards a robust vision system

Modeling issues

- involve a lots of parameters to test i.e. large number of neurons
- are complex and compute-intensive
- are difficult to study and to improve

Potential for streaming architectures

- the algorithms are embarrassingly parallel
- are well-suited due to excellent arithmetic intensity
- exhibit 2D/3D locality

The trick of this trade is **how to leverage their potential for visual perception**, the objective of my thesis.

Visual saliency model

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



Static saliency map



Dynamic saliency map

Visual saliency model

Features

- linearly modeled all the way from the retina to cortical cells
- separation of useful information into two distinct signals that are more efficient to process
- motion compensation estimates and eliminates the camera motion
- motion estimation is used to carry out the motion contrast map

The problem

+ exhibit data parallelism \rightarrow single operation on huge data with no or less dependency - slower performance \rightarrow complex computations involving high

- slower performance \rightarrow complex computations involving high storage and memory usage

↓ Requires real time capability To provide interactivity

Hence, the algorithm is GPU friendly

Related work

- Itti [2002] and Ouerhani [2003] achieved real-time but by simplification of the model
- Longhurst [2006], Peters [2007], Xu [2008] employed only static modalities for the saliency map
- Lee et al. [2007] use both static and dynamic modalities but in a virtual environment

Itti, L.: Real-time high-performance attention focusing in outdoors color video streams. HVEI (2002) Lee, S. et al.: Real-time tracking of visually attended objects in interactive virtual environments. VRST (2007) Longhurst et al.: A guu based saliency map for high-fidelity selective rendering. AFRIGRAPH (2006) Ouerhani, N. et al.: Real-time visual attention on a massively parallel simd architecture. Real-Time Imaging (2003) Peters, C.: Toward 3D selection and skeleton construction by sketching. Eurographics Ireland (2007) Xu, T. et al.: Looking at the surprise: bottom-up attentional control of an active camera system. ICARCV (2008)

Objectives

- To port data-parallel vision algorithm
- To demonstrate the speedup
- To confirm effects of low precision
- To apply different optimizations
- To experience the difficulties
- To incorporate other processes into the model
- To demonstrate the usability of the saliency maps
- To automate the process of porting and optimization

GPU as co-processor

• Stream processing architecture

- Programmability
- Precision
- Power
- New shared memory and synchronization
- Provides texture lookups
- Very high global memory bandwidth
- Accessible, easier to program and manage
- Supported by GPU-specific libraries like CUFFT, CUDPP, CuBLAS, GpuCV
- Already applied in diverse fields





Porting onto GPU Multi-GPU platform

Porting onto GPU

Computer vision algorithms usually comprise of a sequence of filters, hence they are relatively easier to port on to GPU



Introduction Results Conclusion

Porting onto GPU Multi-GPU platform

Porting onto GPU

Computer vision algorithms usually comprise of a sequence of filters, hence they are relatively easier to port on to GPU



- Use texture & shared memory
- Global memory coalescing
- Decrease & optimize use of shared memory
- Substitute math operations
- Remove if and for



Porting onto GPU Multi-GPU platform

Porting onto GPU

Computer vision algorithms usually comprise of a sequence of filters, hence they are relatively easier to port on to GPU



Introduction Results Conclusion

Porting onto GPU Multi-GPU platform

Multi-GPU platform



Introduction Results Conclusion

Porting onto GPU Multi-GPU platform

Multi-GPU platform



Performance gains Effects of lower precision

Conclusion

Speedup



Figure: Speedup for the visual saliency model

Performance gains Effects of lower precision

Conclusion

Profile



Figure: Timings for the two channels of the visual saliency model

Performance gains Effects of lower precision

Conclusion

Precision



Figure: The effect of lower precision support on the result

Universal image index = Q = 99.66%



- Real time capability (\sim 25 fps)
- Created opportunity to extend the model
- Exploited GPUs power without extensive re-structuring
- Without the need for high precision

Activities

International Journal:

Rahman, A.; Houzet, D.; Pellerin, D.; Guyader, N. & Marat, S. Parallel Implementation of a Spatio-temporal Visual Saliency Model. Journal of Real-Time Image Processing, JRTIP. Elsevier (2010)

Conferences:

- Houzet, D.; Huet, S. & Rahman, A. SysCellC: a data-flow programming model on multi-GPU. International Conference on Computational Science, ICCS '10 (2010)
- Rahman, A.; Houzet, D. & Pellerin, D. GPU implementation of motion estimation for visual saliency. Conference on Design and Architectures for Signal and Image Processing, DASIP (2010) {Submitted}

Book chapter:

 Rahman, A.; Houzet, D. & Pellerin, D. Visual Saliency model on Multi-GPU. mei Hwu, W. (ed.) GPU Computing Gems4(Nvidia book). Elsevier (2010)

Presentations:

- Rahman, A. Parallel implementation of visual saliency model on GPU. Presentation during CUDA programming course, Nov, GIPSA-lab (2009)
- Rahman, A. Bio-inspired visual attention model on GPU. Invited speaker at L'utilisation des GPU pour les applications de traitement du signal et des images, GDR-ISIS, 06th May, Telecom Paristech (2010)

Posters:

- Rahman, A.; Houzet, D.; Pellerin, D.; Guyader, N. & Marat, S. Parallel Implementation of a Visual Saliency Model, NVIDIA Research Summit, GTC '09, Oct, San Jose, USA (2009)
- Rahman, A.; Houzet, D.; Pellerin, D.; Guyader, N. & Marat, S. Visual saliency on multi-GPU. CNRS -EEFTIG '10 Ecole d'été francophone de traitement d'image sur GPU, 29 June - 02 July, Grenoble (2010)

Future plan(1)

Porting on to GPU cluster

- to evaluate GPU power consumption
- to measure bandwidth between host and device memories
- to study peak performance single, or multiple CPU-GPU nodes
- to test using different precisions, and various benchmarks
- to evaluate different GPU programming languages and toolkits

Rhone-Alpes CIBLE project in collaboration with LaHC St. Etienne.

Future plan(1)

Porting on to GPU cluster

- to evaluate GPU power consumption
- to measure bandwidth between host and device memories
- to study peak performance single, or multiple CPU-GPU nodes
- to test using different precisions, and various benchmarks
- to evaluate different GPU programming languages and toolkits

Rhone-Alpes CIBLE project in collaboration with LaHC St. Etienne.

Multi-core CPU solution

- to take advantage of the multicore CPUs (OpenMP)
- to support different interconnection networks (OpenMPI)

Future plan(2)

Inclusions in the pathway

- to investigate and model the inclusion of sound
- to contribute to robustness of the exiting model

Future plan(2)

Inclusions in the pathway

- to investigate and model the inclusion of sound
- to contribute to robustness of the exiting model

Creation of a library or tool

- to investigate the correctness of the model for decision making
- to prepare a fast platform used by collaborators for investigating videos

Thank you for your attention!