

GPU Computing Projects

E. Carlinet, J. Chazalon {firstname.lastname@lrde.epita.fr}

Fall 2022

EPITA Research & Development Laboratory (LRDE)



Overview

Objectives

The goals of the project are to:

- apply data-parallelism concepts
- practice with CUDA
- set up a benchmark with a sound evaluation procedure
- present your results in a clear and convincing way

Possible Subjects

Standard Assignment

We propose 1 subject, that most of you should work on:

Implementation of simple object detector in CUDA

More details later in this presentation.

Custom Assignment

For students who are at ease with CUDA, and want to investigate a particular question:

Implementation and performance analysis of *SOME INTERESTING* algorithm in *YOUR PARALLEL PROGRAMMING TECHNOLOGY OF CHOICE*

If you choose this assignment, you must **validate your subject with us.**

Contact us by email ASAP.

Important dates

Deadline	What	Where
Oct 2, 23:59	Group composition (teams of 4)	on Moodle
Oct 9, 23:59	Video submission (see later)	On nextcloud
Oct 23, 23:59	Final project submission (teams of 4)	on Moodle
Oct 24, all day	Oral defense	(either Teams or in presence)

Deliverables

Final Deliverables (1/3)

1. Implementation

- Source code for C++ CPU reference
- Source code for CUDA implementation(s)
- Source code for benchmark tools
- Build scripts (GNU Make, CMake...)

We must be able to reproduce your results

Final Deliverables (2/3)

2. Report

- Description of the problem
 - Detailed if custom subject
 - Quick summary otherwise
- Quick description of the baseline CPU impl.: paper reference, parallel or not, etc.
- Quick description of the baseline GPU impl.: changes from CPU version, **kernels implemented**, etc.
- **Which performance indicators you have used and why**
- **Identification of performance bottlenecks** (with measured indicators, graphs, etc.)
- **For each improvement over the GPU baseline** (implementations):
 - justification of this work regarding performance analysis
 - description of the improvement (e.g. used output privatization instead of global atomics)
 - comparison of the performance with and without this implementation
- **Table with summary** of the benchmark of all variants implemented
- **Summary of the contributions of each team member**

Final Deliverables (3/3)

3. A live lecture / defense

- 15' presentation
- 5' demo
- 5' discussion
- Defenses will be held on Teams (opt. in presence)
- All the team members must be there

We will share with you some images/videos to process during the defense.

Links to each meeting will be shared when all groups are formed.

Submission

Submit code + report + slides on Moodle

<https://moodle.cri.epita.fr/course/view.php?id=948>

Grade Sheet Used for Last Session

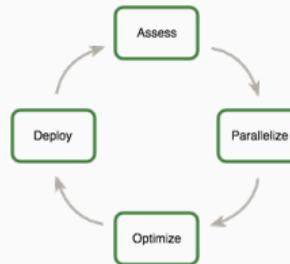
	Travail effectif	Présentation	Note individuelle
Total	10	5	5
Critères	Difficulté technique (bonus/malus) Qualité de l'implémentation (au regard des notions vues en cours) Qualité de l'évaluation, du benchmark (2 implém. CPU, 1+ GPU...) Prise de recul, analyse (identification des goulots d'étranglement...)	Qualité du support Clarté de l'exposé – Pédagogie Mise en contexte Positionnement et état de l'art	Prise de parole Réponse aux questions Implication dans le groupe
Repères			
A [16-20]	Benchmarks cohérents Au moins 3 implémentations (dont 2 variantes GPU) Utilisation de nvprof et/ou nsiight (présentation graphique) Description des techniques utilisées Explication des différences de performance Identification des pistes d'amélioration	Oral de qualité, apprécié de l'auditoire Structure claire Illustrations claires Bonne maîtrise du temps	Maîtrise globale du sujet Réponses pertinentes aux questions Leader de groupe (A+)
B [12-16[Bonne implémentation (versions CPU et GPU avec gain) Au moins un benchmark entre les 2 implémentations Identification d'au moins un facteur de ralentissement (mais pas forcément de solutions) MAIS une seule version GPU	Présentation retenant l'attention (B+) Présentation orale honnête (B-)	Participation honnête au projet Capable de faire appel aux notions de cours
C [8-12[Minimum attendu Version GPU fonctionnelle MAIS manque de prise de recul sur les gains possibles MAIS manque d'analyse de la performance OU pas de version CPU	Présentation passable (C+) Présentation décevante (C-)	Présentation confuse Passif lors des questions Maîtrise incomplète du sujet
D [5-8[Clairement en dessous du travail attendu Implémentation GPU non fonctionnelle	Mauvaise présentation orale	Travail très faible Mauvaise compréhension du sujet
E [0-5[Pas de travail significatif Pas d'implémentation GPU	Inacceptable / absent / ...	Inacceptable / absent / ...

How you should work

Our Expectations

We expect your implementation to be:

- running on GPU;
- correct, i.e. it produces an acceptable result.

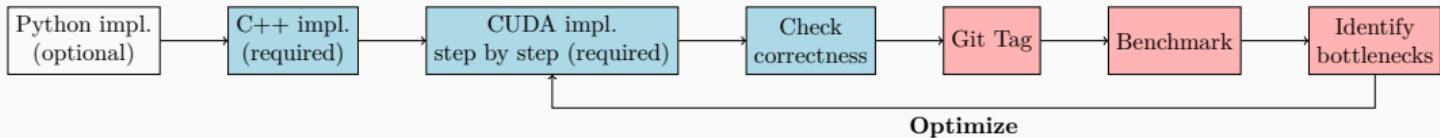


Do not try to make it fast at first, just make it work.

Then, try to apply NVidia's Assess, Parallelize, Optimize, Deploy (APOD) design cycle as described in the CUDA C++ Best Practices Guide:

1. identify the part of the code which is responsible for the bulk of the execution time;
2. get a parallel version of the code (assumed to be sequential at first);
3. optimize the performance of the parallel code;
4. measure the performance of the new code.

Recommended approach



Some hints:

- Have a working (*slow*) C++ reference implementation first (and keep it forever)
- Tag (*git tag*) the versions of your program before any optimization (useful to track and benchmark ideas)
- Try optimizations step by step so that you can tell which ones are the most important

Project Outline for Standard Performance Analysis

Broad Outline	Concrete Example
Choose an application	Mandelbrot
Determine the most time-consuming part of the app	Global atomics
Determine one or more data-parallel approaches to solve the problem	Tiling...
Create multiple implementations of the approach	One naive version, one version with shared memory...
Benchmark the implementations	Record memory transfer time, kernel time, utilization, FLOPS, etc.
Relate results to course concepts	Identify the cause of the bottleneck (memory or compute bounding)

Object Detection

Objective

Object detection

Reference image	After	Algorithm output (boxes)
		

In practice: at least a working CUDA implementation which can process 1 image

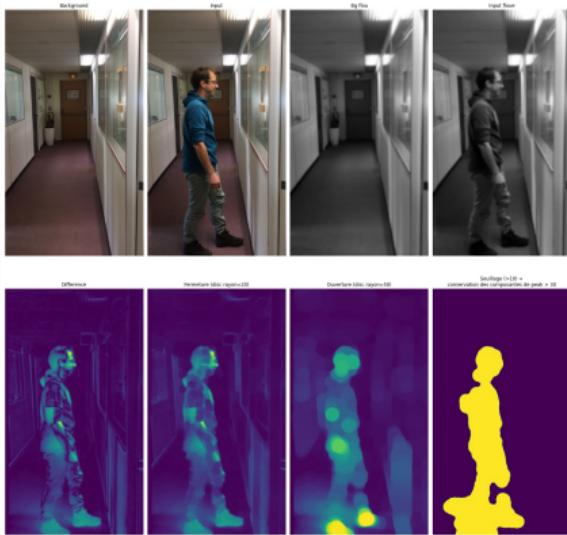
Inputs and Outputs

Your program should be able to take an image (or an image sequence) and a reference image as inputs, and return the coordinates of the boxes ([x y width height]) of the detected objects for each image.

```
$ ./mybin reference.jpg input-0001.jpg input-0002.jpg input-0003.jpg
{
    "input-0001.jpg" : [
        [0, 0, 10, 10],
        [15, 15, 42, 42]
    ],
    "input-0002.jpg" : [],
    "input-0003.jpg" : [
        [0, 0, 10, 10],
        [15, 15, 42, 42],
        [51, 42, 69, 99]
    ]
}
```

You can use a python script to display the results in a nice way (e.g. a video with rectangles as overlays).

Algorithm steps



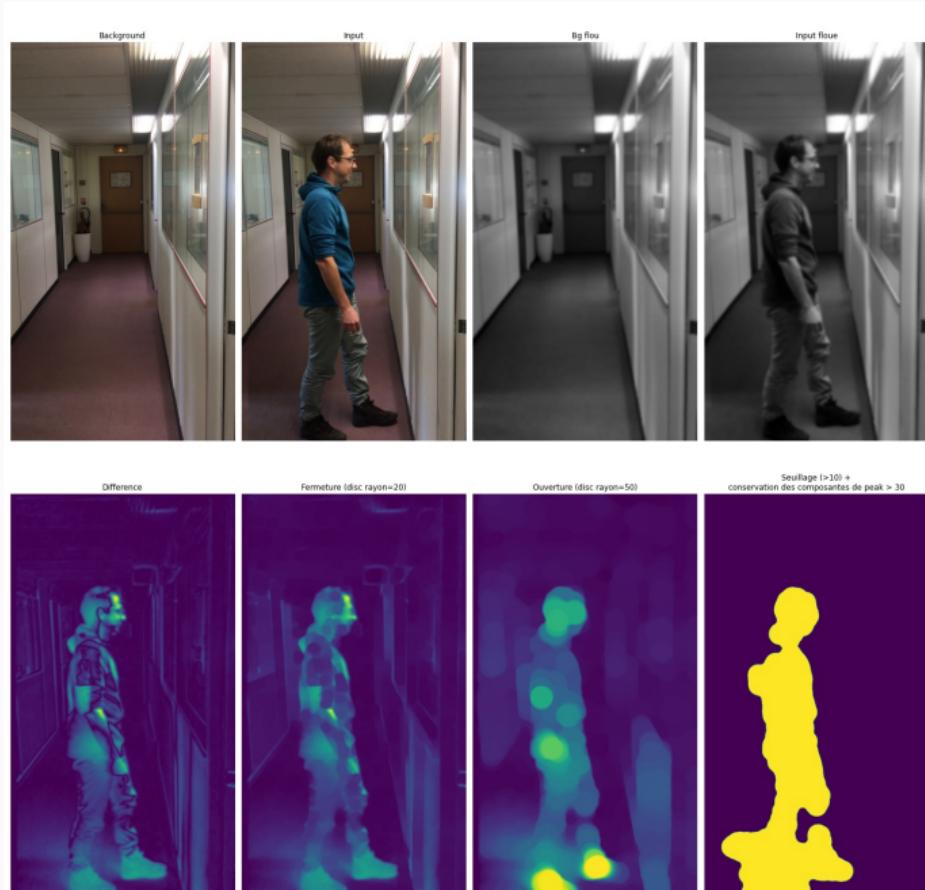
1. Convert the image in grayscale 
2. Smooth the background image g_0 and the new image g_1 (e.g. with a Gaussian filter¹) 
3. Compute the difference between images
$$d = |g_0 - g_1|$$
 
4. Perform morphological closing/opening with a disk (or rectangle) to remove non-meaningful objects² 
5. Threshold the image and keep only the blobs (connected components) with high peaks ³
6. Output the bounding boxes 
7. Display the image with a red border around the detected blobs (on CPU with Python)

It is up to you to find which kernels you need, and which pattern they correspond to.

¹ A Survey of Gaussian Convolution Algorithms, Pascal Getreuer, <https://www.ipol.im/pub/art/2013/87/article.pdf>

<https://canvas.colorado.edu/courses/61439/pages/morphological-opening-and-closing> ³ Mathematical morphology and its applications to image processing, Soille, 2001

Zoom on the steps



Connected component labeling + peaks

```
init(L); // L[p] ← p
while L has changed:
    propagate(L); // L[p] = min L[n] (n neighbor of p)
    relabel(L); // Make label continuous
    A ← get_peaks_and_bbox(L);
```

Expected work

To get the average grade:

- all deliverables (code, report, slides)
- working CPU version
- working GPU version
- benchmark their respective speed and compute relevant indicators (occupancy, L1/L2 cache hit rates...). The **FPS** rate should be reported.

To get a better grade:

- implement more variants
- perform more analysis
- use cuda streams

Dataset

Collaborative dataset for videos:

<https://cloud.lrde.epita.fr/s/oE5t9TpweK3cDBj>

Each **group** has to upload 1 video of ~15s (use your phone) **stabilized** where an object/person is introduced in the scene (the first frame will be used as the reference image).