

# Development of USTH interactive virtual world using Unreal Engine

**DO Duy Huy Hoang** Supervisor **TRAN Giang Son** *University of Science and Technology*

2019-06-18

## List of Figures

1	Execution time of Stena on 1080p images . . . . .	6
---	---	---

# Contents

<b>Introduction</b>	<b>1</b>
Context and Motivation . . . . .	1
Objectives . . . . .	1
<b>Technical details</b>	<b>1</b>
Utilities functions . . . . .	1
Convolution . . . . .	2
Sorting . . . . .	2
Output . . . . .	2
<b>Implementation GPU</b>	<b>3</b>
Methods . . . . .	3
Data formatting . . . . .	3
Kernel call configuration . . . . .	3
Loading data into shared memory . . . . .	4
Convolution and output . . . . .	4
Other steps . . . . .	5
Benchmark results . . . . .	5
Execution time comparison . . . . .	5

# Introduction

## Context and Motivation

### Objectives

We denote Steganographie as **Stena**, for short.

In this project, we have done the following:

- Stena with CPU
- Stena with GPU and CUDA (5-15x faster) using shared memory
- Hiding multiple characters in an image
- Can use filter any size (with width, height  $\leq 31$ , odd numbers)
- Automated testing for correctness and speed
  - Benchmark at the end

## Technical details

In this part, we highlight the important ideas and implementation details in the program. Note that this section is about CPU version, the GPU version will be explained later.

Overall, Stena is divided into 4 main phases:

- Initialization: We allocate memory for the arrays we use during computation
- Convolution: the filter is convoluted with the image to create image V. Pixel of V at the border (that cannot fit the filter) will have value 0
- Sorting:  $8*n$  pixels with the largest values in V is selected, where n is the length of the input string. Each pixel will contain 1 bit of that string
- Output: given  $8*n$  positions, the characters are hidden in the original image using the described method

The convolution section is the most important part in this project. The final two is self-explanatory in the source code.

## Utilities functions

- We use ppm\_lib.h for image reading/write. Since the memory does not include any memory management, we have to allocate/free memory manually
- For measuring time, we use C++ chrono libraries. The utility class is implemented in my\_timer.h
- Definitions of important variables:
  - H, W: height and width of the image
  - M, N: height and width of the filter
  - n: length of the hidden string

- $K = 8*n$ : the number of pixels needed to hide the string, 1 bit per pixel

## Convolution

The CPU implementation is straightforward. We loop over each pixel  $(i,j)$  and place the center of the filter on top of it. Then, we loop over each cell of the filter to compute the convolution sum. That sum will be the new value of  $V(i,j)$

Note that for pixels at the border, the filter cannot fit inside the image. In that situation,  $V(i,j) = 0$ .

Time complexity:  $O(H*W*M*N)$

## Sorting

In this part, we need to find  $8*n$  pixels  $(i,j)$  with the largest  $V(i,j)$ .

The simplest way to this is to create an array size  $H*W$ , where each element contains a value  $V(i,j)$  and its index in array  $V$  (which is implemented as a 1D array, like `PPMImage`). After that, we sort the array decreasingly and choose the first  $8*n$  pixel position to hide the string.

Time complexity for sorting is:  $O(H*W*\log(H*W))$

This is wasteful because  $n \ll H*W$ . For example, for a 1080p images, there are  $1920*1080 = 2073600$  pixels to sort. If we want to hide a 10-character string, we only need the largest 80 pixels. Sorting the entire array is a waste.

Therefore, we use QuickSelect algorithm to partially sort and find only the  $K$  largest elements. This is a popular and common algorithm, and more details can be seen in the code's comments.

//TODO: add more code desb hereeee

Time complexity:  $\sim O(H*W*\log(K))$

## Output

This section is straightforward. For each character in the string, we hide each of its bit in the Least-significant Bit of the selected pixels, starting from the pixel with the largest  $V(i,j)$ . The method is exactly the same as the project's description.

Using bit-wise operations, it is simple to manipulate the bits of a number. More details can be found in the code's comments, inside function `hideString()` and `getString()` of file `stena_cpu.cpp`.

//TODO: add more code desb hereeee.

# Implementation GPU

## Methods

We only use GPU acceleration for convolution. For sorting and outputting, we use the same functions from the CPU version. Since convolution accounts for the overwhelming majority of the execution time, this improvement still speeds up the program greatly.

We used Tiled 2D Convolution, which uses Shared memory to improve speed. The idea is from this lecture. In this project, we assume the filter sizes satisfy:  $M, N \leq 31$

We define 2 new variables:

- $\text{rowpb} = \text{TILE\_DIM} - M + 1$ : the number of rows a block process
- $\text{colpb} = \text{TILE\_DIM} - N + 1$ : the number of columns a block process

## Data formatting

In the `ppm_lib.h` library, a Pixel is represented by a struct:

```
struct PPMPixel {
    unsigned char red, green, blue;
}
```

However, in the GPU version, we use raw array of bytes. That means 3 consecutive elements in the array represent 1 pixel. For comparison, in an array of pixels, we do the following to access the color channels of the  $i$ -th pixel:

- PPMPixel: `a[i].red`, `a[i].green`, `a[i].blue`
- Raw data: `a[3*i]`, `a[3*i+1]`, `a[3*i+2]`

We do this for simplicity, and because raw data is easier to handle in kernel code.

## Kernel call configuration

Our kernel launch blocks of  $32 \times 32$  threads, `TILE_DIM=32` in the code, each block processes a tile of size  $\text{rowpb} \times \text{colpb}$  pixels. Each block processes  $\text{colpb}$  continuous columns. After it finishes with the first  $\text{rowpb}$  rows, it moves on to the next rows, and so on. Therefore, we need to launch enough blocks to cover all columns, which equals to:  $\text{roundup}(W / \text{colpb})$ . The kernel launch looks like this:

```
int columnsPerBlock = TILE_DIM - filtW + 1;
dim3 grid((imgW + columnsPerBlock - 1) / columnsPerBlock, 1, 1);
dim3 block(TILE_DIM, TILE_DIM, 1);
myconv2dCuda <<< grid, block, 2 * filtH * filtW * sizeof(float) >>>
    (imgH, imgW, gdata, filtH, filtW, gfilter, gV);

// todo, desc more details hereee
```

### Loading data into shared memory

In Tiled 2D convolution, first we define a 2D array for shared memory, called *smem*.

```
__shared__ float smem[TILE_DIM][TILE_DIM];
```

Then, we need to define some variables (more details in the code comment). // TODO describe code

```
// rowsPerBlock = TILE_DIM - filtH + 1; Formula: roundup(a/b) = (a + b - 1) / b.
const int loop = (imgH + (TILE_DIM - filtH + 1) - 1) / (TILE_DIM - filtH + 1),
              stride = TILE_DIM - filtH + 1,
              halfH = filtH / 2, halfW = filtW / 2;
const int mycol = blockIdx.x * (TILE_DIM - filtW + 1) + tidy,
          mycolOut = mycol + halfW;
int myrow = tidy, myrowOut = tidy + halfH;
```

Recall that at each step, a block process rowpb rows of array V. So, loop is the number of steps needed to compute H rows of the image. For each block, the first thing to do is loading data from global memory (the image) into shared memory. Since we launch blocks of 32×32 threads, which match smem exactly, loading data is very simple.

```
// load image data to shared memory
if (mycol >= imgW || myrow >= imgH) smem[tidy][tidy] = 0;
else {
    int pixelPos = 3 * cell(myrow, mycol, imgW);
    smem[tidy][tidy] = pixels[pixelPos] + pixels[pixelPos + 1]; // red, green
}
__syncthreads();
```

The variables myrow, mycol represent the current pixel position of a thread. The **top-left** corner of the filter is placed on this position (unlike the center of the filter in the CPU version). Therefore, if it's outside of the image, its value is 0. Otherwise, we read its value from the global memory to shared memory. **Note:** cell(i,j,w) = i\*w + j. Given the row and column of a pixel, this macro turns it into the actual address of the pixel in the data array.

### Convolution and output

For each thread, the top-left corner of the filter is placed at (myrow, mycol). Therefore, the output pixel is placed at (myrowOut, mycolOut), which is shown in the second image of the previous subsection.

```
// convolute. Note that all threads in a warp access the same filter[cell(i,j,f
// So, we use constant memory for better speed.
if (tidy < TILE_DIM - filtH + 1 && tidy < TILE_DIM - filtW + 1 &&
    // the top-left corner of the filter is put here, and it must fit inside th
```

```

        myrowOut < imgH - halfH && mycolOut < imgW - halfW) // It must output to a
    {
        float tmp = 0;
        for (int i = 0; i < filtH; i++)
            for (int j = 0; j < filtW; j++)
                tmp += smem[tidx + i][tidy + j] * filter[cell(i, j, filtW)];

        V[cell(myrowOut, mycolOut, imgW)] = tmp;
    }

    // update indexes
    myrow += stride;
    myrowOut += stride;

```

Finally, convolution is very simple. First, we only consider threads that can fit the filter inside the image and outputs to a pixel inside the image. Then, each thread loop over the filter and multiply it with the corresponding pixel (in shared memory). We noticed that all threads in a warp always access the same element of the filter at each step in this loop, so the filter is put in **constant memory** for better performance. By doing this, instead of reading the data serially, the GPU can issue a **broadcast**, meaning all threads get the required data immediately. Finally, the corresponding pixel  $V(\text{myrowOut}, \text{mycolOut})$  is updated, and blocks move on to process the next set of rows.

### Other steps

After computing  $V$ , the sorting and outputting sections use the same CPU functions as the CPU Stena version.

## Benchmark results

For benchmarking, we use FullHD images (1920x1080) because this is a very common resolution. We measure the execution time of Stena at different filter radius (for benchmarking, we use square filters).

The benchmark is done on a computer with 8750H CPU at 3.1GHz, 16GB dual-channel 2666MHz DDR4, 1060 6gb, and Windows 10. The program is compiled in Visual Studio 2019 Release mode (all optimizations on).

### Execution time comparison

You can easily see that the CPU version's running time has quadratic growth with respect to the filter radius. Meanwhile, the GPU version is very fast in most cases.



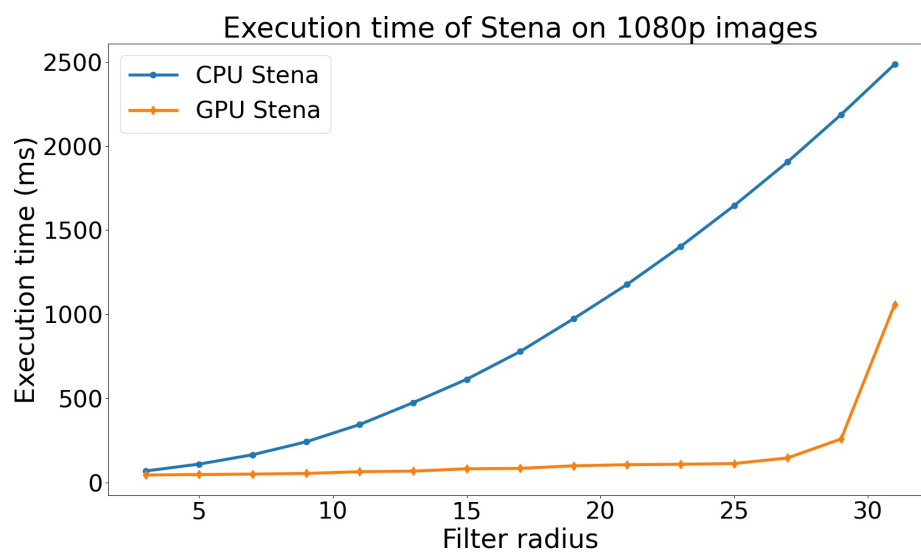


Figure 1: Execution time of Stena on 1080p images