# Labwork 2: Hello, CUDA!

Phi Doan Minh Luong - 2440046

September 2025

## 1 Preprocessing

We will use the Lenna image for this labwork.



We load the image using the OpenCV library, then flatten the image into a 1D array of RGB values using the Numpy library

```
1  rgb = cv2.imread("Lenna.png")
2  imageHeight, imageWidth, _ = rgb.shape
3  print(imageHeight, imageWidth)
4  pixelCount = imageHeight * imageWidth
5  rgb = np.reshape(rgb, (pixelCount, 3))
6  print(rgb.shape)
```

## 2 CPU

We implement grayscale using the CPU

```
1  def grayscale_cpu(src, dst):
2    for i in range(pixelCount):
3      g = np.uint8((int(src[i, 0]) + int(src[i, 1]) + int(src[i, 2])) / 3)
4      dst[i, 0] = dst[i, 1] = dst[i, 2] = g
5    return dst
```

The average processing time after running 10 times is 0.4067821502685547 seconds. The resulting image is:



# 3 GPU

First, we defined a function to convert an image to grayscale on the GPU.

```
1  @cuda.jit
2  def grayscale_gpu(src, dst):
3    tidx = cuda.threadIdx.x + cuda.blockIdx.x * cuda.blockDim.x
4    g = np.uint8((src[tidx, 0] + src[tidx, 1] + src[tidx, 2]) / 3)
5    dst[tidx, 0] = dst[tidx, 1] = dst[tidx, 2] = g
```

Then, we follow the flow of the CUDA program, calculate the average time processed after running each block size 10 times

```
1  devOutput = cuda.device_array((pixelCount, 3), np.uint8)
2  devInput = cuda.to_device(rgb)
3  blockSize = [32, 64, 128, 256, 512, 1024]
4  timeProcess = []
5
6  for b in blockSize:
7    gridSize = math.ceil(pixelCount / b)
8    avgtime = []
9    for _ in range(10):
```

```
10        start = time.time()
11        grayscale_gpu[gridSize, b](devInput, devOutput)
12        avgtime.append(time.time() - start)
13        hostOutput = devOutput.copy_to_host()
14        grayscale_gpu_image = hostOutput.reshape((imageHeight, imageWidth, 3))
15    cv2.imwrite(f"Lenna_gray_gpu_{b}.png", grayscale_gpu_image)
16    timeProcess.append(np.mean(avgtime))
17    print(f"Time-processed-with-block-size-{b}:-{np.mean(avgtime)}s")
```
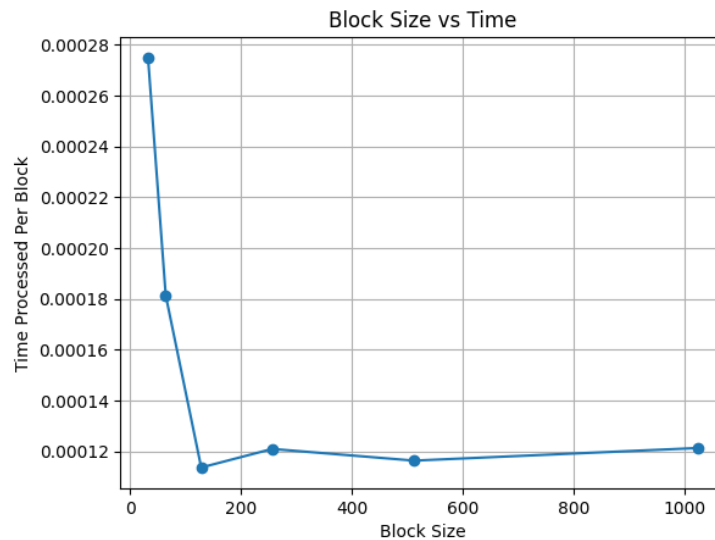
The resulting images are all the same. The average time results are:

```
1  Time processed with block size 32: 0.0002749919891357422s
2  Time processed with block size 64: 0.00018134117126464844s
3  Time processed with block size 128: 0.00011363029479980469s
4  Time processed with block size 256: 0.00012094974517822266s
5  Time processed with block size 512: 0.0001163482666015625s
6  Time processed with block size 1024: 0.000121307373046875s
```

Here, I plot a graph of block size vs time



At the small block sizes (32 and 64), the processing time per block is higher. This is because small block sizes may not fully utilize the GPU, leading to a higher time per block.

At the optimal block sizes (128 to 512), it used threads and cores more efficiently, reducing the processing time.

At the high block size (1024), there is a slight increase in processing time.