

Acceleration of JPEG Decoding Process using CUDA

Rushikesh Tade

Department of Electronics and Telecommunication
Pune University
Pune - 411 007

Saniya Ansari

Department of Electronics and Telecommunication
Pune University
Pune - 411 007

ABSTRACT

In this paper we have implemented efficient JPEG (Joint Photographic experts group) Decoder on GPU (Graphic Processing Unit) using NVIDIA CUDA (Compute Unified Device Architecture) Technology. This decoder is capable of decoding images of Ultra HD resolution with superfast speed GPU is used to assist the CPU for time consuming tasks. In this paper IDCT module which consumes 70 to 80 percent of computation time is implemented on GPU. An asynchronous parallel execution between the CPU and the GPU is used at a same time to improve the JPEG decoder acceleration rate. In this work, the JPEG decoder based on the CUDA performs decompression of images of size 2560 x 1600 pixels and below. Finally the results are shown with respect to different sized images and consumed time for decoding. The results show that this decoder faster in multiple times than the decoder in CPU.

General Terms:

GPUGPU, CUDA, JPEG, DCT

Keywords:

CUDA, JPEG, RLE ,HUFFMAN Decoding

1. INTRODUCTION

In recent year many with the development of new technology data storage cost is also increased. The more familiar example would be the screen resolution of the TV sets. We can see how the popularity of high resolution TV sets is been increasing over the years. Since the data is increased and so do the quality of material. This has resulted in more data to be stored and transferred via communication media. Which has lead us to the various compression techniques for the data to be transferred and stored. Now a days we can find many standards for the image compression. JPEG (Joint Photographic Experts Group) is one them which has very popular for photographic image decoding [1]. In computing, JPEG is a common method of lossy compression for digital images, especially for those images produced by digital photography. Adjustment in the compression can be done in JPEG, allowing a selectable trade-off between storage size and image quality. JPEG usually achieves 10:1 compression with little perceptible loss in quality of image.

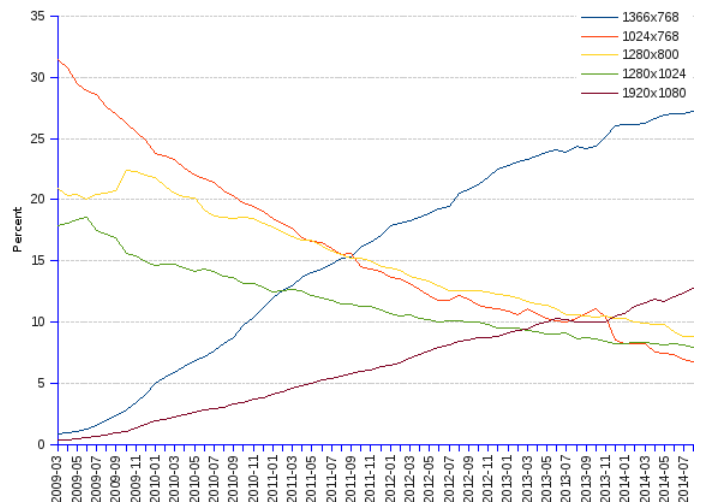


Fig. 1. Percent share of screen resolution in Europe market [3]

Since its very popular in image compression many software and hardware techniques have been developed to improvise the speed of JPEG decoding process. To the hardware resources for the JPEG decoder, Some common efficient ways can be obtained using the ,FPGA (Field Programmable Gate Array) or other ASIC (Application Specific Integrated Circuit) resources, DSP (Digital Signal Processor) [2].With increasing computations speed of GPU along with the CUDA framework a software decoder can be efficiently implemented.

2. CUDA TECHNOLOGY

General-Purpose computing on Graphics Processing Units (GPUGPU) is referred to techniques where calculations traditionally done by the Central Processing Unit (CPU), handed over to the Graphics Processing Unit (GPU). Earlier, the GPU was used only to accelerate certain parts of the graphics pipeline, but now it can reduce the CPU load and/or increase the processing throughput for general purpose scientific and engineering computing. A GPU acceleration model is shown in figure 3. GPU-accelerated computing gives good application performance by giving load of compute-intensive portions of the application to GPU, while the rest of the code still runs on the CPU.

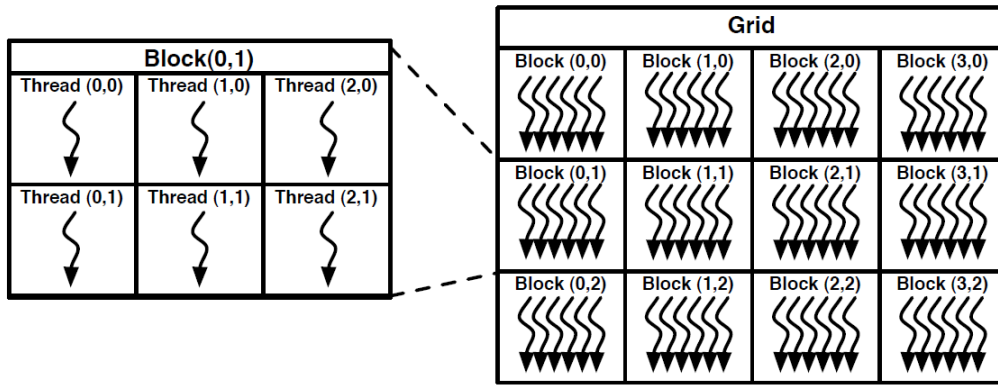


Fig. 2. CUDA execution model

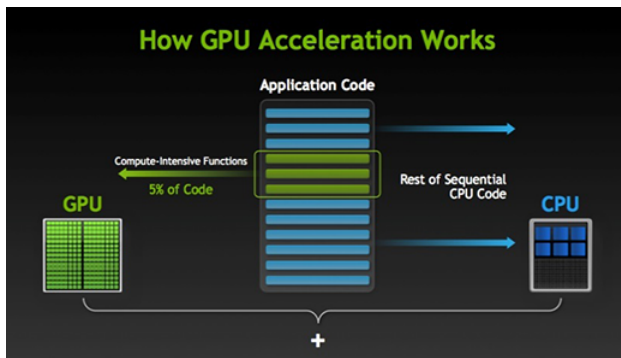


Fig. 3. GPU acceleration Model [4]

From a user's view, applications simply run significantly faster. In order to ease the usage of GPUs for programmers not familiar with the graphics pipeline, the CUDA language was created by NVIDIA. CUDA is a programming interface to parallel architecture for general purpose computing. It is an extension to the C language, with a programming model easily understood by programmers already familiar with threaded applications. This interface is a set of library functions which can be coded as extension of the C language. A compiler produces executable code for CUDA device. The CPU sees a CUDA device as multi-core co-processor. The CUDA design does not have memory restrictions of GPGPU. One has an access of all memory available on the device using CUDA with no restriction on its representation though the access times vary for different types of memory [5]. The GPU works with a very high amount of threads which run simultaneously. Threads are run in different batches called thread blocks [6]. And the CUDA kernel function is essentially based on the block as a unit. [7] In the same block, there are maximally 512 threads. As the threads of the same block are executed into the same streaming multiprocessors (SM), they can share data with each other by shared memory. The number of active blocks is not more than 8 in each SM. Generally, the number of active threads is not more than 768. There are six types of memory in the CUDA programming model: Global Memory, Shared Memory, Local Memory, Register, Constant Memory and Texture Memory [4]. The global memory is located in the memory and it can be both read and write by devices, but it has large memory access latency. However constant memory

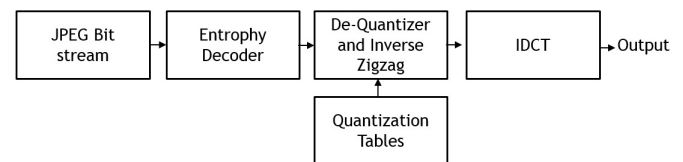


Fig. 4. JPEG decoding process on CUDA

and texture memory are cached, the access speed is faster than global memory. But they are only read by devices. The registers are on the GPU caches, execution units of GPU can access to register with very low latency, the register is the basic unit of the register file and each register file only has 32 bit. When the program has many private variables, it can use the local memory. When the program needs data communication in a block, the shared memory is a good choice, it is also on cache and the access speed of shared memory almost as fast as the register. So properly using the shared memory to reduce other kinds of memory access is a good way to reduce delay caused by memory access. [8]

2.1 nVidia Performance Primitives (NPP) Library

The NVIDIA Performance Primitives library (NPP) is a collection of GPU-accelerated video, image, and signal processing functions that deliver up to 5 to 10 times faster performance than comparable CPU-only implementations. Using this library, we can take advantage of over 1900 image processing and approximately 600 signal processing primitives to achieve significant improvements in the performance of application. [9] We have used this library to boost our performance.

3. JPEG DECODING PROCESS AND IMPLEMENTATION ON CUDA

JPEG decoder usually consists of an entropy decoder which usually consists of Huffman decoder. Before the start of the actual image in the decoder's stream of input data, the encoder has placed the tables it used for encoding the image. This means that the decoder can first extract the required tables from its input data stream, in order to process the rest of the input data with the help of these tables. The first step is to reverse the entropy encoding process,

which produces the quantized DCT coefficients, with the help of the quantization tables, the DCT coefficients can be de-quantized and finally be transformed back via the IDCT process. Figure 4 shows the decoding process of the JPEG. [10] [11]

The Huffman decoder includes the decoding of the run-length encoding (RLE), the entropy coding of the DC(Direct Current) coefficients and AC(Alternating Current) coefficients. At the same time, in order to decode the DC coefficients, the current DC variable is to add the former DC variable of a color component unit, and it involves a large number of logical operations. So the Huffman decoding is more suitable to use the CPU to complete [12] [13]. In this paper, the de-quantizer and the inverse zigzag operation are done with the Huffman decoder module. The de-quantizer is accomplished by multiplying the value with the quantization table when the coefficients is finished by the Huffman decoder model. Then the inverse zigzag is implemented with the help of look up table. Since the hardware configuration of the GPU is not suitable for logic operations. The GPU will bring large overhead in the implementation of the control flow instruction. At the same time, this model is also so simply done in the Huffman model, so this module is more efficiently done by the CPU than by the GPU with the CUDA. The whole 2-D IDCT have more data computation in the JPEG decoder than any other block. The 2-D IDCT in the decoded image is no correlation between each sub-graph, so which makes it computable independently. Therefore, there are different levels of parallelism in the step of the IDCT transform. So we can use the CUDA to finish it. In this work have used CUDA NPP library for the processing of the data in CUDA is been used. This gives users advantage for the processing of the most of the data into the GPU memory. Which has resulted as boost of performance as discussed in the Results. Flowchart for the whole process is shown in figure 5. First a bit stream is processed in the CPU and different markers are detected from it. Using this bit stream we have separated out the data such as Huffman tables, Quantization tables and data bits. After that we have performed De-quantization and Huffman decoding of the data. First, we use the texture memory to store the Huffman decoder data for the input of IDCT module. With this way we can reduce the access latency comparing to global memory. Then we move with the constant memory to store some constant data such as cos matrix and it makes the constant data access latency as fast as the register. Secondly, in the same block, we use the shared memory to store the output data from the row 1-D IDCTs. Since there are eight sub IDCTs executing in the GPU, there are 512 intermediate data for the next 1-D IDCT step. As a result, we put the 512 intermediate data into the shared memory and it will also reduce the access latency of data. Which has resulted into producing the DCT data. Preceding that we have used IDCT module for which outputted the pixel values of the data.

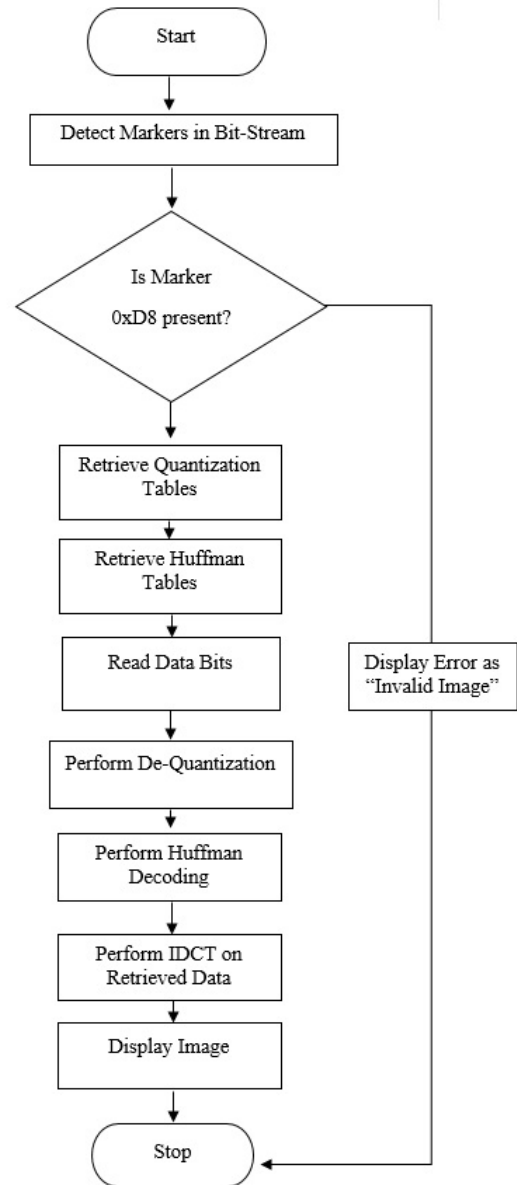


Fig. 5. JPEG decoding process on CUDA

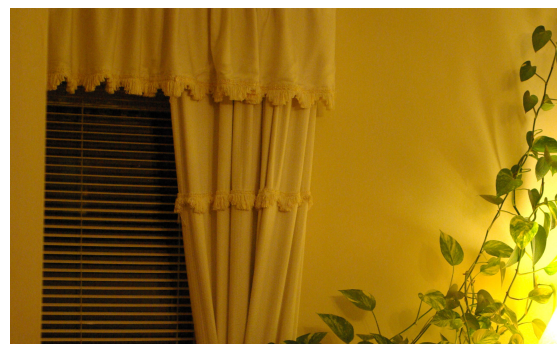


Fig. 6. Input Image

4. EXPERIMENTAL SETUP AND RESULTS

We have used images of three size for the decoding process. Input image is shown in 6. The three resolutions which are used in for experimentation are 600 x 522, 1920 x 1080, and 3240 x 2160. The Results are calculated with by calculating the time required for the execution a module. We have also compared the results with the method proposed by Ke Yan [14]. Ke Yan has used two methods for the implementation of the CUDA decoder it. Synchronous and Asynchronous execution of the CUDA decoder.

Table 1. Results for the execution of Image Decoder

Image Size	CPU	GPU asynch	GPU synch	NPP
600 x 522	31 ms	16 ms	21 ms	11.72 ms
1920 x 1080	187 ms	89 ms	117 ms	44.85 ms
3240 x 2160	621 ms	279 ms	373 ms	124 ms

Table 2. Results for the execution of IDCT Module

Image Size	CPU	GPU	NPP
600 x 522	18 ms	0.38 ms	0.28 ms
1920 x 1080	125 ms	2.54 ms	3.33 ms
3240 x 2160	422 ms	8.54 ms	5.55 ms

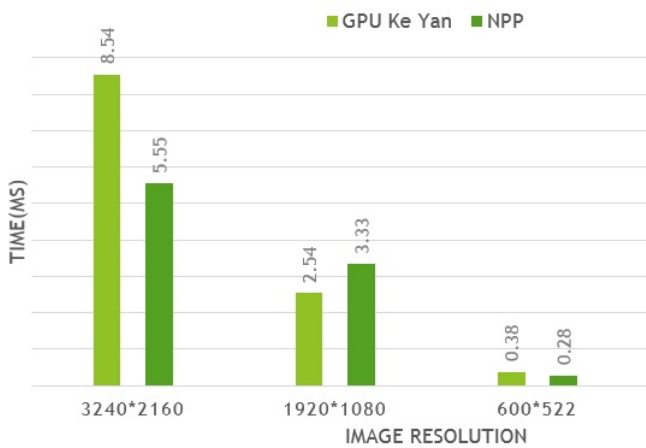


Fig. 7. IDCT Results

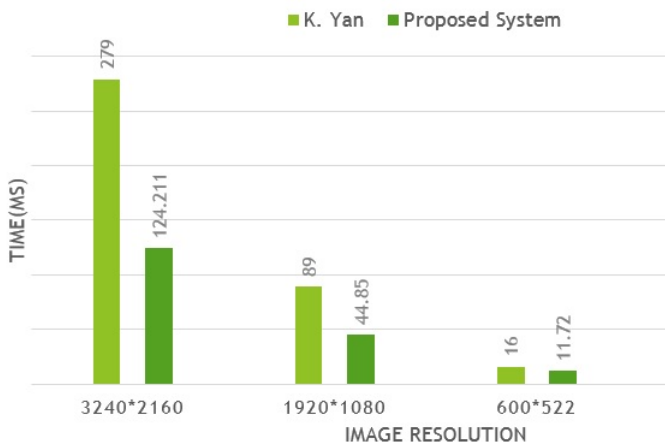


Fig. 8. Results for the Decoder

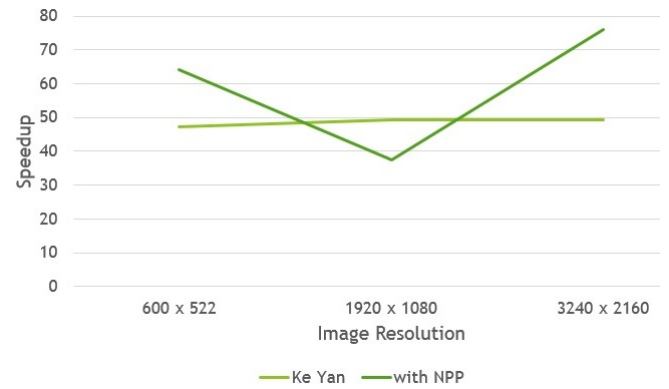


Fig. 9. IDCT Results

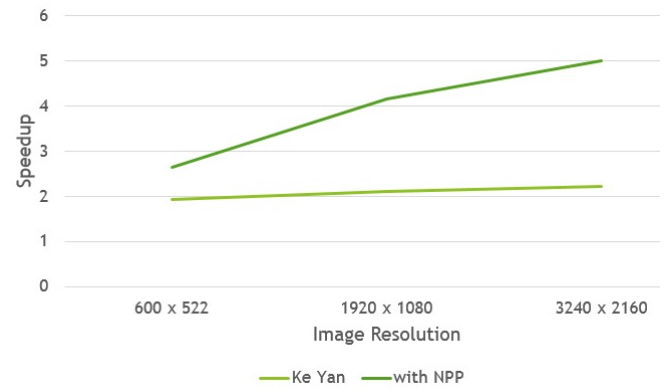


Fig. 10. Results for the Decoder

These results show that with this method we get around 26x more speed-up than the previous method and also around 3x more speedup in decoding process. This is mainly because most of our operations are performed in GPU memory.

5. CONCLUSION

In this paper, we introduced the CUDA technology and analysed the basic JPEG decoder models. By taking advantage of the CUDA technology and NPP libraries, we can speed up the JPEG decoding by much more than speed than the traditional CUDA based JPEG Decoders. The experimental results show that the CUDA-based JPEG decoder can save about 70% time than the CPU-based realization and the model of the IDCT decoder realized on the GPU can be 49 times faster than the CPU-based realization.

Our future work will be targeted in the direction of improving the JPEG decoder and we will move towards optimization of JPEG resizing process which will be achieved by doing the implementation of both decoder and encoder unit on CUDA.

6. ACKNOWLEDGEMENT

The authors would like to thank the Department of Electronics and Telecommunication Engineering of Dr. D.Y. Patil School of Engineering, as well as researchers for making their resources available and teachers for their guidance. We are thankful to the authorities Board of Studies Electronics and Telecommunication Engineering of Savitribai Phule Pune University. We are also thankful to reviewer for their valuable feedback and suggestions. We also thank the college authorities for providing the required infrastructure and support. Finally, we would like to extend a heartfelt gratitude to friends and family members.

7. REFERENCES

- [1] Kun-Bin Lee and Chi-Cheng Ju. A memory-efficient progressive jpeg decoder. In *VLSI Design, Automation and Test, 2007. VLSI-DAT 2007. International Symposium on*, pages 1–4. IEEE, 2007.
- [2] Jahanzeb Ahmad, Kamran Raza, Mansoor Ebrahim, and Umar Talha. Fpga based implementation of baseline jpeg decoder. In *Proceedings of the 7th International Conference on Frontiers of Information Technology*, page 29. ACM, 2009.
- [3] Areppim AG, stats of screen resolution eu. http://stats.areppim.com/stats/stats_screenresxtime_eu.htm. Accessed: 2015-04-12.
- [4] NVIDIA, gpu accelerated computing. <http://www.nvidia.com/object/what-is-gpu-computing.html>. Accessed: 2015-04-25.
- [5] Pawan Harish and PJ Narayanan. Accelerating large graph algorithms on the gpu using cuda. In *High performance computing–HiPC 2007*, pages 197–208. Springer, 2007.
- [6] CUDA Nvidia. Compute unified device architecture programming guide. 2007.
- [7] Hong Biao Li. A new efficient method for dct8x8 with cuda. In *Applied Mechanics and Materials*, volume 681, pages 231–234. Trans Tech Publ, 2014.
- [8] Shane Cook. *CUDA programming: a developer's guide to parallel computing with GPUs*. Newnes, 2012.
- [9] NVIDIA, nvidia performance primitives. <https://developer.nvidia.com/NPP>. Accessed: 2015-04-25.
- [10] Jingqi Ao, Sunanda Mitra, and Brian Nutter. Fast and efficient lossless image compression based on cuda parallel wavelet tree encoding. In *Image Analysis and Interpretation (SSIAI), 2014 IEEE Southwest Symposium on*, pages 21–24. IEEE, 2014.
- [11] KS Priyadarshini, GS Sharvani, and SB Prapulla. A survey on parallel computing of image compression algorithms jpeg and fractal image compression. *IJITR*, pages 78–83, 2015.
- [12] Jeong-Woo Lee, Bumho Kim, Jungsoo Lee, and Ki-Song Yoon. Gpu-based jpeg2000 decoding scheme for digital cinema. In *Advanced Communication Technology (ICACT), 2014 16th International Conference on*, pages 601–604. IEEE, 2014.
- [13] Bart Pieters, Charles-Frederik Hollemeersch, Jan De Cock, Peter Lambert, and Rik Van de Walle. Data-parallel intra decoding for block-based image and video coding on massively parallel architectures. *Signal Processing: Image Communication*, 27(3):220–237, 2012.
- [14] Ke Yan, Junming Shan, and Eryan Yang. Cuda-based acceleration of the jpeg decoder. In *Natural Computation (ICNC), 2013 Ninth International Conference on*, pages 1319–1323. IEEE, 2013.