

# Parallel Video Processing

[Parallel Programming Final Project Report] \*

Po-Han Chen<sup>†</sup>  
National Chiao Tung University  
j.lin013@gmail.com

Cheng Sun<sup>‡</sup>  
National Chiao Tung University  
s2821d3721@gmail.com

Yu-Wen Pwu<sup>§</sup>  
National Chiao Tung University  
yuwen41200@gmail.com

## ABSTRACT

Video processing typically requires immense amounts of computational resources. Devices supporting up to 4K resolution are currently hitting the market by storm. A greater demand for computation has thus arisen. Therefore, we set out to explore ways to boost the efficiency of video processing through a series of experiments. These experiments were conducted with a wide assortment of tools(or APIs) on multiple platforms. We've managed to achieve what we consider maximum possible speedups on our test platforms.

## 1. INTRODUCTION

Initially, we were interested in whether parallelizing simple image processing tasks could deliver any speedup. A number of image processing tasks rely on matrix multiplications which take only constant number of computations per pixel. This gives  $O(N)$  complexity and likely higher memory access frequencies.

Seeing these computations are reasonably simple, we suspected that compilers might optimize programs to an extent where speeding up with parallelization might not be possible. To test this out, we designed and conducted 2 experiments, *lightup* and *lightup2*. The former does simple arithmetic calculations on each pixel, while the latter multiplies each pixel by  $1.125^{(x+y)/128}$ . From these experiments, we concluded that increasing the amount of computations per thread does raise speedups and that different platforms(or OSes) along with different optimization options are not the reasons for the subtle speedups in *lightup*.

Our research into parallel video processing branched out from here. We decided to have a closer look at the feasibility of using FPGAs to speedup video processing. We consider FPGAs a suitable choice for heterogeneous computing be-

cause FPGAs are essentially massively parallel processors. Furthermore, fully customized hardware design can lead to maximized optimization. Past research has also showed that using FPGAs as accelerators can realize better performance than using GPGPUs. In the end, we were able to obtain nearly 1000x speedup on a 1920\*1080 dataset. This was also 100x faster than other approaches.

Unsatisfied with *lightup* experiments, we've found and tested a tool called *pmbw* - Parallel Memory Bandwidth Benchmark. It's a tool written in assembly that tests important memory characteristics including maximum bandwidth on our test platforms. After examining the test results from *pmbw*, we believe a speedup of up to 3.1x should be attainable.

We then turned our attention to real-world problems. We decided to use White balance as our video processing application. The White balance algorithm we used takes  $O(N)$  time and was simple enough for us to experiment with. We used 5 different parallelization approaches: Pthread\_TDM, OpenMP\_TDM, CUDA, CUDA\_TDM and TaskParallel. We will detail these approaches in later contexts. While CUDA still delivers poor results due to frequent global memory accesses, Pthread and OpenMP managed to yield 3x speedups on the main test platform and 2x speedups on the secondary test platform. Both platforms confirmed the test results from *pmbw*. We concluded that certain hardware characteristics do pose an "invisible ceiling" especially in I/O bound applications and CUDA on consumer-grade GPUs might not be suitable for this type of applications.

## 2. PROPOSED SOLUTIONS: LIGHTUP SERIES (CPU/GPU)

### 2.1 *lightup*

For each pixel,  $[R', G', B'] = 2 * [R, G, B] + 5$   
Then normalize the RGB values back to  $[0, 255]$ .  
This would increase the overall color contrast.

### 2.2 *lightup2*

For each pixel,  $[R', G', B'] = [R, G, B] * 1.125^{(x+y)/128}$   
Then normalize the RGB values back to  $[0, 255]$ .  
The images would become mostly 0's (ie. darken) except for the bottom-right corner.

\* We will discuss our parallelization methods with White balance altogether.

\*  
†  
‡  
§

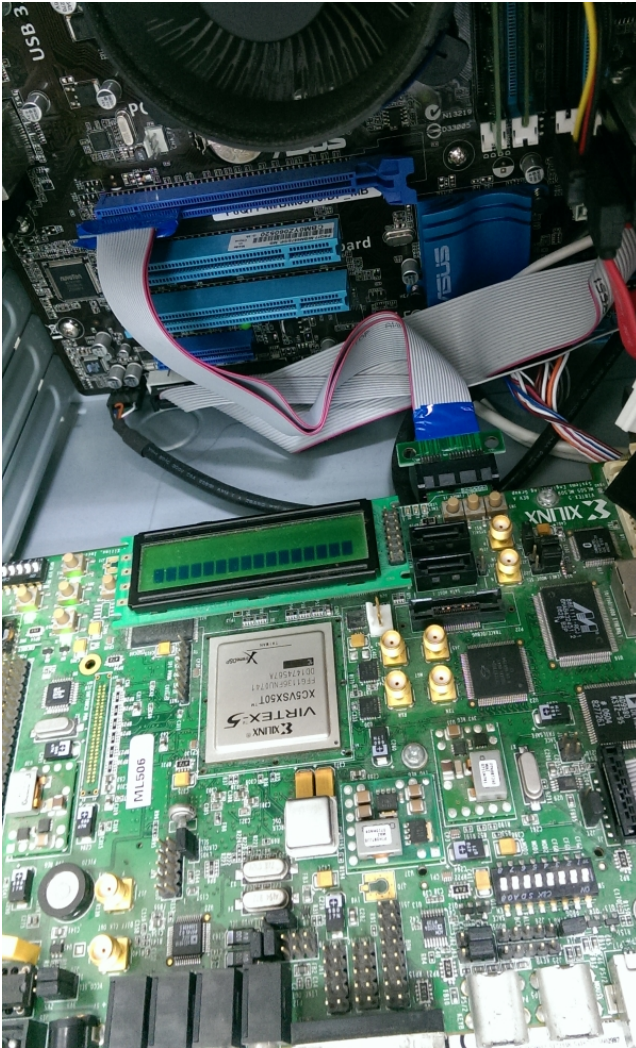
### 3. PROPOSED SOLUTIONS: LIGHTUP2.5 (FPGA)

#### 3.1 lightup2.5

For each pixel,  $[R', G', B'] = [R, G, B] * 1.125^{1024}$

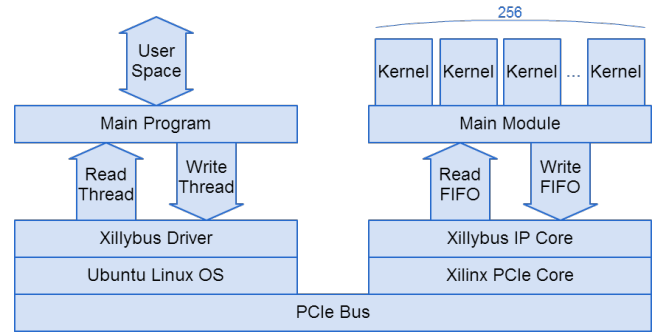
#### 3.2 Implementation

We employ RTL (register-transfer level) design by the Verilog HDL (hardware description language), and we use the Xillybus IP core [2] for data transmission over the PCIe interface. Our targeted board is Xilinx Virtex-5 ML506 Evaluation Platform (w/ XC5VSX50T FPGA). We use C++ for software design. The host programs run on 64-bit Linux distributions.



#### 3.3 Computation Platform

The PC-end (host) is on the left side, whereas the FPGA-end (device) is on the right side.



##### 3.3.1 Data Interface

To hosts, the device looks just like a file. So we can use low-level (POSIX) file I/O functions to read/write data from/to the device, i.e. the FPGA board.

##### 3.3.2 Multithreading on Software

For best performance, there are also 2 threads implemented on the host. The first is used to send data, the other is used to receive data.

#### 3.4 Kernel Instantiation

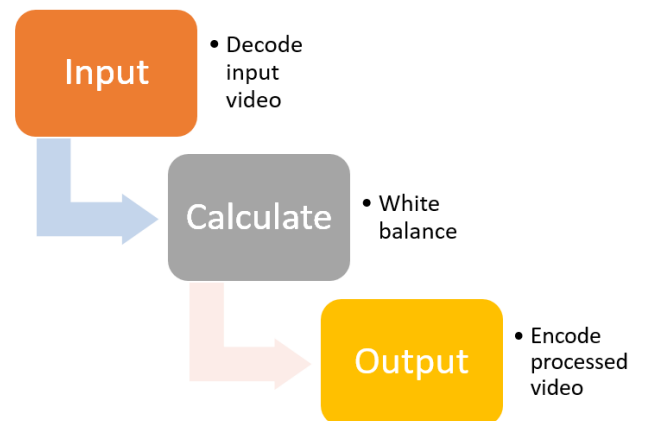
We use a *generate block* to instantiate 256 kernel modules, namely, 256 simultaneously running threads. Each kernel module calculates a single color value, so the FPGA can calculate at most 256 values at the same time.

A Finite-State Machine was used. In it, there are 4 states: IDLE\_STATE, RECV\_STATE, EXEC\_STATE, and SEND\_STATE. After the program starts, the circuit goes to the second state. The second state waits and receives 256 color values, until then, it goes to the third state. The third state processes the 256 received values in parallel, then it goes to the last state. The last state waits and sends 256 new values back, and goes to the first state again.

### 4. PROPOSED SOLUTIONS: WHITE BALANCE

#### 4.1 Block diagram

Input and Calculate are parallelizable.

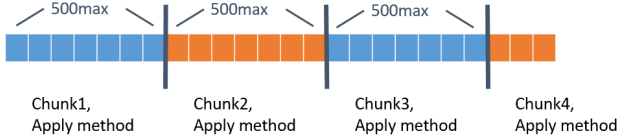


## 4.2 Algorithm

1. Calculate the averages of Red, Green and Blue values, as denoted by (AvgR, AvgG, AvgB)
2. Using Green as reference, for each pixel,  
 $Red_{new} = Red * AvgG / AvgR$   
 $Blue_{new} = Blue * AvgG / AvgB$
3. If  $Red, Blue_{new} > 255$  Then  $Red, Blue_{new} = 255$

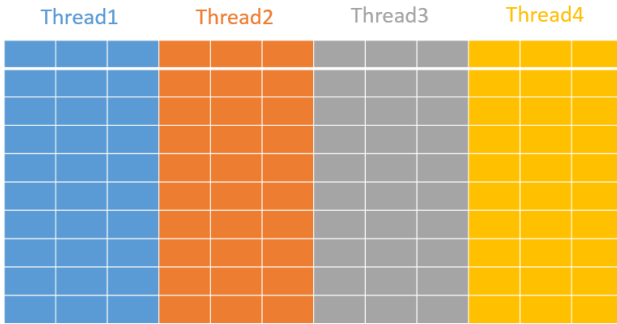
## 4.3 Initialization

With scalability in mind and to prevent memory exhaustion, we divide the input video into chunks with each containing no more than 500 frames.



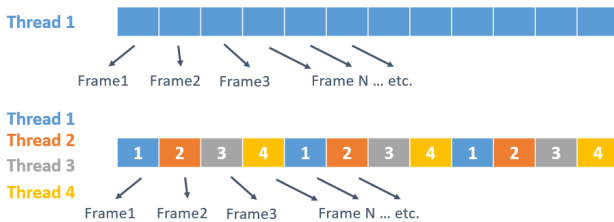
## 4.4 Parallel method1: Pthread

We process one frame at a time. We let each thread process a portion of columns.



## 4.5 Parallel method2: Pthread\_TDM

Utilizing C++11 threading libraries, we divide the input video by time. Each thread is assigned a collection of video frames as indicated by the following figure.

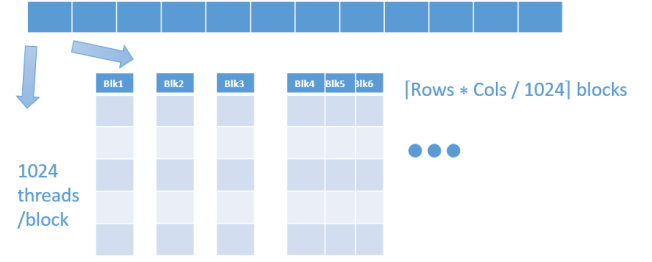


## 4.6 Parallel method3: OpenMP\_TDM

We also divide the input video by time, but with OpenMP as the API.

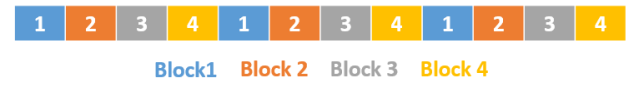
## 4.7 Parallel method4: CUDA

We process one frame in each iteration (as in memcpy operation pairs). Each thread processes one pixel.



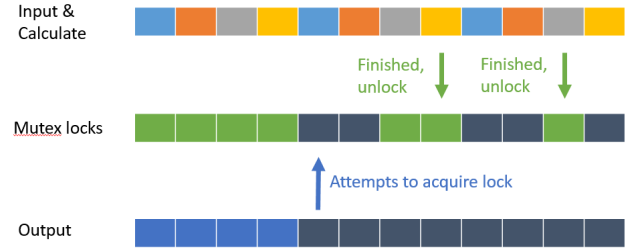
## 4.8 Parallel method5: CUDA\_TDM

Very similar to previous approaches, we divide the input video by time where each thread block processes one frame.



## 4.9 Parallel method6: Pthread\_with\_TaskParallel

Extended from Pthread\_TDM, rather than waiting for Input and Calculate to finish, we create a separate thread specifically for output.



## 5. EXPERIMENTAL METHODOLOGY

### 5.1 Environment specifications

APIs: OpenCV 2.4.11, OpenCV 3.0

FPGA: ML506 Evaluation Platform

OS	CPU	RAM	Compiler
Windows 10	i7-3770(4C8T)	16GB	VS12, gcc 5.2.0
Ubuntu 14.04	i5-4200H(2C4T)	8GB	gcc 4.8.4
OSX El Capitan	i5-4260U(2C4T)	8GB	gcc 5.2.0
Ubuntu 15.10	2 vCores	2GB	gcc 5.2.0

GPU (on i7 PC): GeForce GTX 670, CUDA v7.5, CUDA Capability 3.0

### 5.2 Input datasets

#### 5.2.1 Lightup & Lightup2:

360\*240 / 1280\*720 / 1920\*1080, 148 frames (MP4).

### 5.2.2 *Lightup2.5:*

256\*256 / 1920\*1080\*3 (Pure data).

### 5.2.3 *White balance:*

1280\*720, 1422Frames (MP4, AVI).

## 5.3 Output

### 5.3.1 *Lightup & Lightup2:*

During visual inspection, the programs display output video in realtime.

### 5.3.2 *Lightup2.5:*

Pure data (arrays) of sizes 256\*256 / 1920\*1080\*3.

### 5.3.3 *White balance:*

If output is opted in (through C macro definition), the programs produce videos in the same resolution as input in AVI format.

## 6. EXPERIMENTAL RESULTS

### 6.1 *lightup*

Slowdown when input is small. Tiny speedup even when input is large. Similar results across OpenCV 2.4, Ubuntu, OSX.

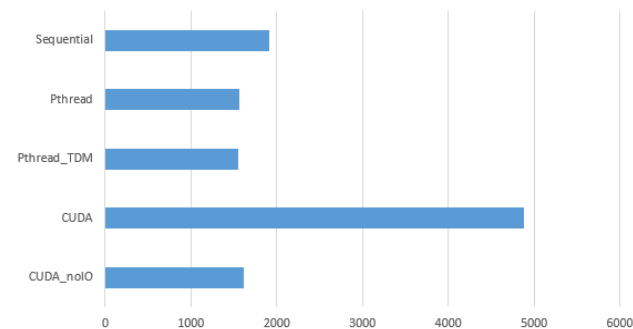


Figure 1: Lightup (Windows 10, OpenCV3, VS12, 4 Threads)

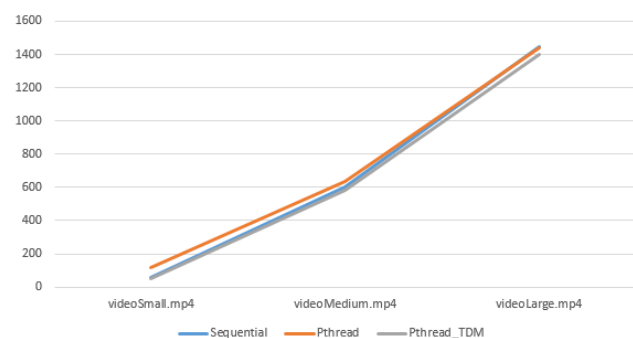


Figure 2: Lightup (Windows 10, OpenCV3, Cygwin-gcc 5.2.0, 4 Threads)

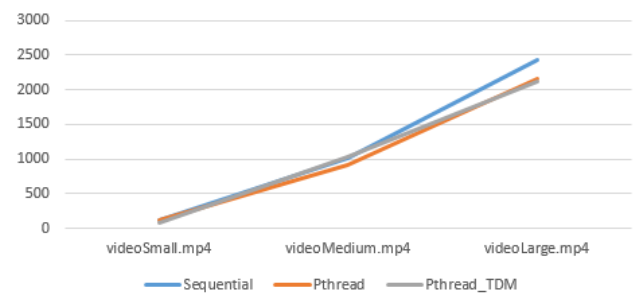


Figure 3: Lightup (OSX El Capitan, OpenCV3, gcc 5.2.0, 4 Threads)

### 6.2 *lightup2*

Results are clearly better. Time Division seems to be better than Frame Division.

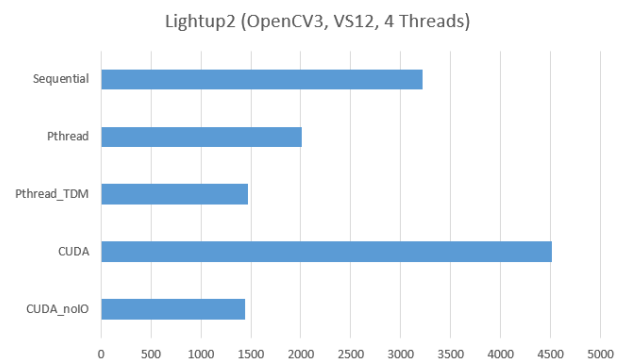


Figure 4: Lightup2 (Windows 10, OpenCV3, VS12, 4 Threads)

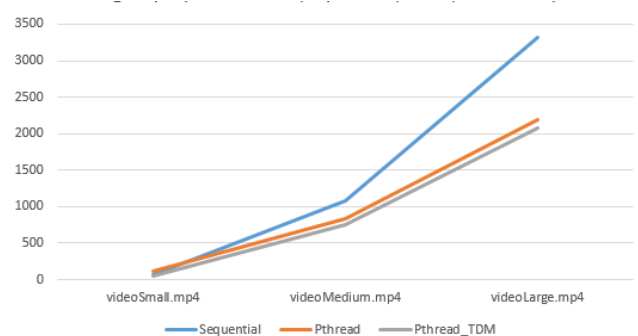


Figure 5: Lightup2 (Windows 10, OpenCV3, Cygwin-gcc 5.2.0, 4 Threads)

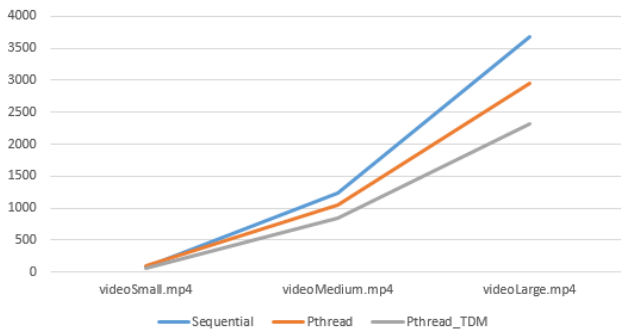
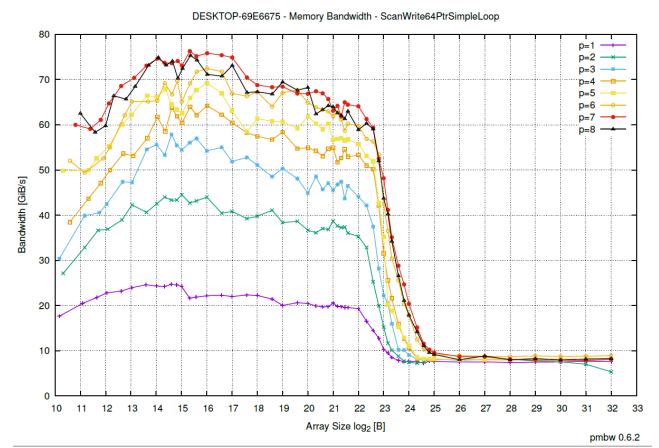


Figure 6: Lightup2 (Ubuntu 14.04, OpenCV3, gcc 4.8.4, 4 Threads)



### 6.3 lightup2.5 (FPGA)

A whopping, near 1000x speedup was observed.

Table 1: lightup2.5 Runtimes

Size	Sequential	OpenMP	CUDA	FPGA
Small	179.096ms	67.667ms	891.333ms	<b>0.534ms</b>
Large	15873.200ms	4678.667ms	1799.333ms	<b>16.104ms</b>

Table 2: lightup2.5 Speedups

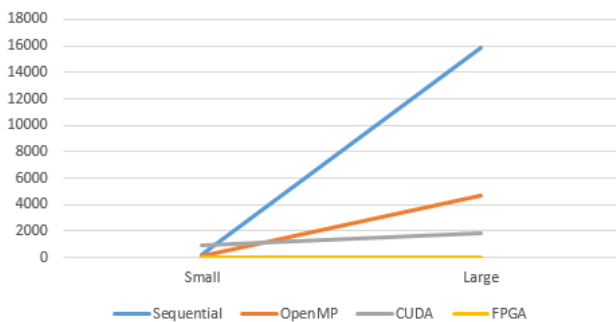
Data Size	Sequential	OpenMP	CUDA	FPGA
Small	1.00x	2.65x	0.20x	<b>335.39x</b>
Large	1.00x	3.39x	8.82x	<b>985.67x</b>

### 6.5 White balance



Figure 7: White balance Results: Before and After

Lightup 2.5



### 6.4 pmbw

From the results below, we conclude that I/O is parallelizable and speeding up is definitely possible. We think the upper bound for speedup for I/O bound applications is roughly 3.1x on the main test platform.

2 threads:  
Pthread\_TDM: 1.69x speedup  
OpenMP\_TDM: 1.77x speedup

White Balance (i7-3770, 2 Threads, VS12, Win10)

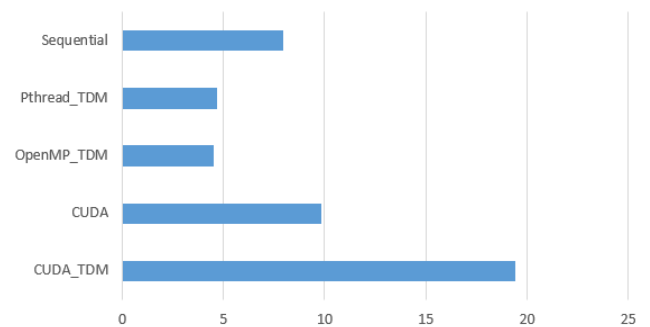


Figure 8: White balance with 2 threads

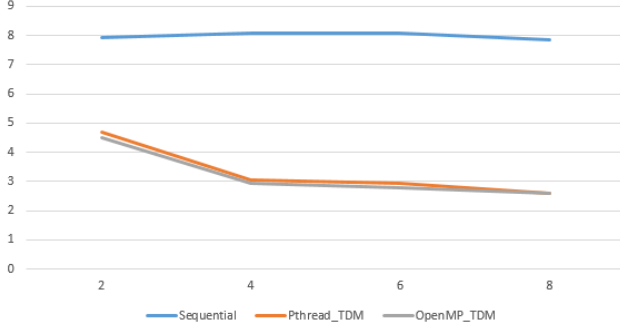
### 6.6 White balance and more threads

We can see our approaches toward white balance do benefit from more threads.

The most speedups we can get correspond to the results from pmbw.

**Table 3: White balance Speedups with more threads**

Threads	1	2	4	6	8
Pthread_TDM	1.00x	1.69x	2.49x	2.76x	3.02x
OpenMP_TDM	1.00x	1.77x	2.76x	2.92x	3.00x



**Figure 9: White Balance (i7-3770, 2 to 8 Threads, VS12, Win10)**

## 6.7 White balance and other platforms

Positive results across different OSes and compilers.

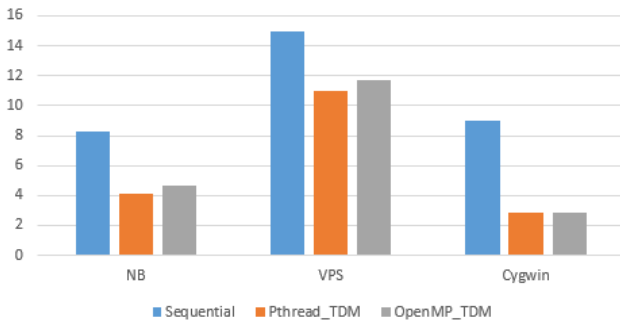
NB: Ubuntu 14.04, i5-4200H(2C4T), 8GB, gcc 4.8.4

VPS: Ubuntu 15.10, 2 vCores, 2GB, gcc 5.2.0

Cygwin: Windows 10, i7-3770(4C8T), 16GB, Cygwin-gcc 5.2.0

**Table 4: White balance Speedups on other platforms**

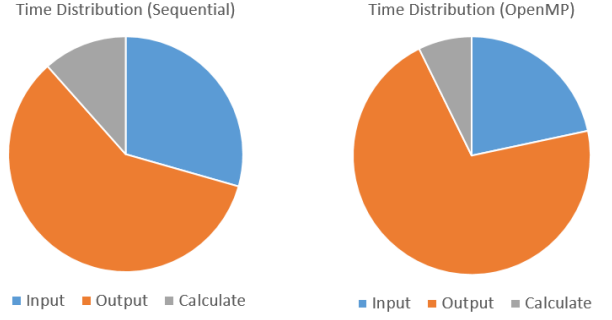
Platform	Seq	NB	VPS	Cygwin
Pthread_TDM	1.00x	2.00x	1.36x	3.13x
OpenMP_TDM	1.00x	1.76x	1.27x	3.20x



**Figure 10: White Balance (NB, VPS, Cygwin)**

## 6.8 White balance Profiling results

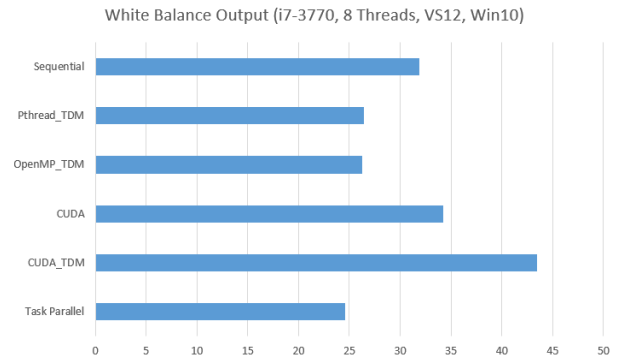
When output is taken into account, the whole process is output-bound.



**Figure 11: Profiling results**

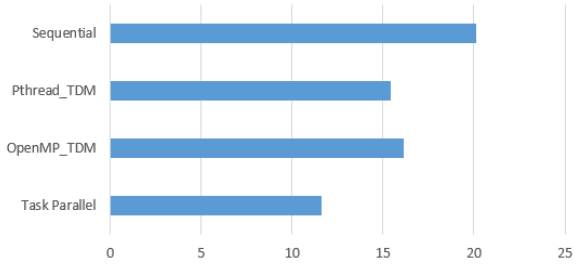
## 6.9 White balance when output is accounted

Pthread\_with\_TaskParallel delivers better results(ie. higher efficiency) when output is taken in account.



**Figure 12: White balance with task parallel on main PC**

**White Balance Output (i5-4260H, gcc 4.8.4, Ubuntu 14.04)**



**Figure 13: White balance with task parallel on NB**

## 7. RELATED WORK

- Ching-Chih Weng, Homer Chen, and Chiou-Shann Fuh. *A Novel Automatic White Balance Method For Digital Still Cameras*. 2005 IEEE International Symposium on Circuits and Systems.
- Dennis Lin, Xiaohuang (Victor) Huang, Quang Nguyen, Joshua Blackburn, Christopher Rodrigues, Thomas Huang,



Minh N. Do, Sanjay J. Patel, and Wen-Mei W. Hwu. *The Parallelization of Video Processing*. IEEE Signal Processing Magazine, November 2009.

[2] Xillybus Ltd. (2016). *An FPGA IP Core for Easy DMA over PCIe with Windows and Linux [Online]*. Available: <http://xillybus.com>

- Mihalis Psarakis, Aggelos Pikrakis, Giannis Dendrinou. *FPGA-based Acceleration for Tracking Audio Effects in Movies*. 2012 IEEE 20th International Symposium on Field-Programmable Custom Computing Machines.
- Matthew Jacobsen, Yoav Freund, Ryan Kastner. *RIFFA: A Reusable Integration Framework for FPGA Accelerators*. 2012 IEEE 20th International Symposium on Field-Programmable Custom Computing Machines.

## 8. CONCLUSIONS

Our experimental results show that certain hardware characteristics, especially memory bandwidth, are closely related and relevant to the efficiency of programs with linear complexities. We experienced a maximum 3.1x speedup on the Win10 desktop and a maximum 2x speedup on the Ubuntu laptop. Both confirmed the pmbw test results from each of their platforms. It seems true that speedups do have upper bounds which depends on the hardware.

Additionally, we think we should keep simple calculations on the CPU end. Using NVIDIA's Nsight profiling tool, we were able to identify huge memory bottlenecks over the course of the program execution. It is likely futile to parallelize simple linear-complexity algorithms as the time it takes to copy the data between host and device would be relatively significant. Things could be better on a workstation-grade GPU though.

Moreover, it seems that FPGAs handle parallel programming extremely well. It certainly wasn't easy to develop on FPGAs, but the staggering results made it all worthwhile. It's probably safe to say that FPGAs are great tools for applications demanding massive calculations.

Lastly, the video codec used for input also seems to matter. When the input video is encoded in mp4 format, parallelizing the input procedure actually slows down the whole process. It took us quite a bit of effort to accomplish the results we have right now. A reasonable guess for this would be, perhaps non-sequential reads caused slight delays in between and longer decoding times potentially triggered unwanted queueing effects.

## 9. CREDITS

- We would like to thank Prof. Yi-Ping You, Assistant Professor at the Department of Computer Science, National Chiao Tung University and the teaching assistants for an excellent course (Parallel Programming, Fall 2015). The knowledge we gained throughout the course helped us tremendously.
- We would also like to thank Prof. Chun-Jen Tsai, Associate Professor at the Department of Computer Science, National Chiao Tung University for lending the FPGA board (ML506 Evaluation Platform) to us and for all the advices he offered during the research.

## 10. REFERENCES

- [1] Jason Su. *Illuminant Estimation: Gray World*. <http://web.stanford.edu/~suja-son/ColorBalancing/grayworld.html>