

EE 382C: Multicore Computing

Parallel GPU based Algorithms for Image Processing

Wenbo Xu, Wenwen Zhang, Yichi Zhang

I. ABSTRACT

II. INTRODUCTION

III. GAUSSIAN FILTER

IV. OPTIMIZATION

A. *Pageable vs. Pinned Memory*

Host data allocations are pageable by default, which means can be paged in/out between RAM and disk. However, GPU cannot access data directly from pageable memory, but from pinned memory, which means page-locked. Hence, whenever a data transfer is invoked on pageable memory, the CUDA driver has to allocate a temporary pinned memory array to copy host data and then transfer it to the device.

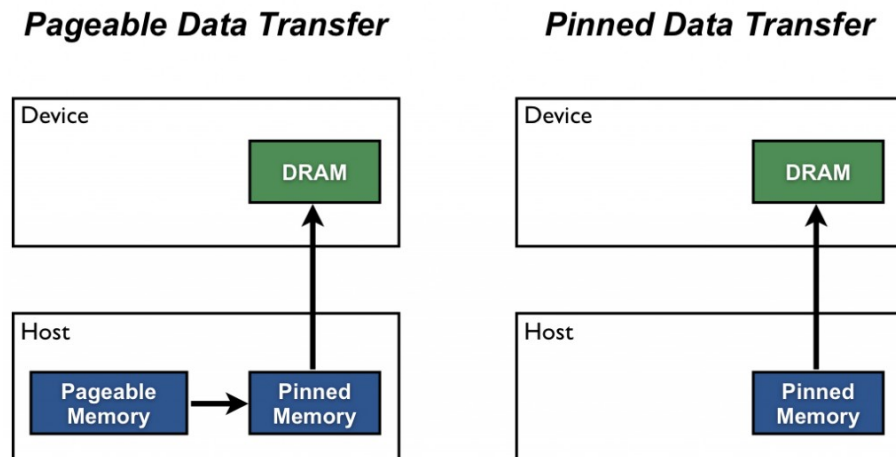


Fig. 1: CUDA data transfer

We can avoid the cost of this overhead by using pinned memory for host instead of pageable memory. In this case, we use *cudaMallocHost()* and *cudaFreeHost()*. Compare to the *malloc()* and *free()*, *cudaMallocHost()* and *cudaFreeHost()* are more expensive with additional overheads. Then, the question has been raised about how should we made the tradeoff. According to figure below, pinned memory is faster when the size of data to be transfered is larger than 16MB.

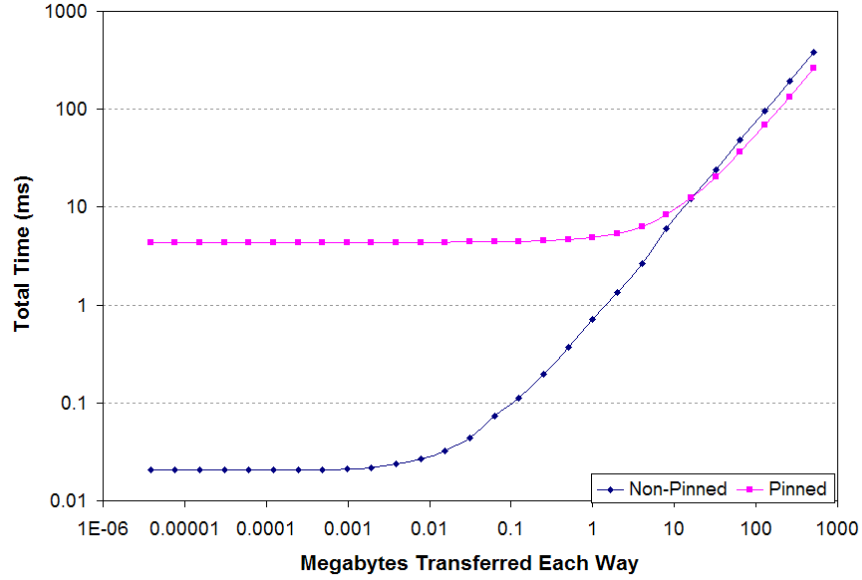


Fig. 2: Time required to allocate, transfer to the GPU, transfer back to the CPU, and deallocate pinned and non-pinned memory.

[2]

This doesn't mean we should never use pinned memory when the amount of data to be transferred is less than 16MB. One example is the asynchronous memory copy, *cudaMemcpyAsync()* can be used only with pinned memory. The details of how asynchronous memory copy would be used to improve the efficiency will be discussed in next section.

B. Streams

A stream is defined as a sequence of operations in that execute in issue-order on the GPU. CUDA operations, which are kernel operations and memory operations, in same streams are ordered and in different streams can overlap. By default, all operations are in default stream. The following code is used to specify which stream the operation is in.

```

for (int i = 0; i < nStreams; ++i) {
    int offset = i * streamSize;
    cudaMemcpyAsync(&d_a[offset], &a[offset], streamBytes, cudaMemcpyHostToDevice, stream[i])
}

for (int i = 0; i < nStreams; ++i) {
    int offset = i * streamSize;
    kernel<<<streamSize/blockSize, blockSize, 0, stream[i]>>>(d_a, offset);
}

for (int i = 0; i < nStreams; ++i) {
    int offset = i * streamSize;
    cudaMemcpyAsync(&a[offset], &d_a[offset], streamBytes, cudaMemcpyDeviceToHost, stream[i])
}

```

The code shown above is an example of invocation of multiple streams in CUDA. The code is divided into three parts, memory copy from Host to Device, kernel invocation, and memory copy from Device to Host. For memory copies, *cudaMemcpyAsync()* is used. As described in last section, we have to allocate pinned memory using *cudaMallocHost()*. This method places the transfer into the stream and returns immediately. It is up to the device to schedule streams when the corresponding resources are free. It allows us to put memory transfer operations and kernel operations into the stream the same time and allow them to run concurrently.[3]

REFERENCES

- [1] Harries, M. (2012, December). How to Optimize Data Transfers in CUDA C/C++. Retrieved from <https://devblogs.nvidia.com/parallelforall/how-optimize-data-transfers-cuda-cc/>
- [2] Boyer, M. Choosing Between Pinned and Non-Pinned Memory. Retrieved from <https://www.cs.virginia.edu/~mwb7w/cuda'support/pinned'tradeoff.html>
- [3] Harries, M. (2012, December). How to Overlap Data Transfers in CUDA C/C++. Retrieved from <https://devblogs.nvidia.com/parallelforall/how-overlap-data-transfers-cuda-cc/>

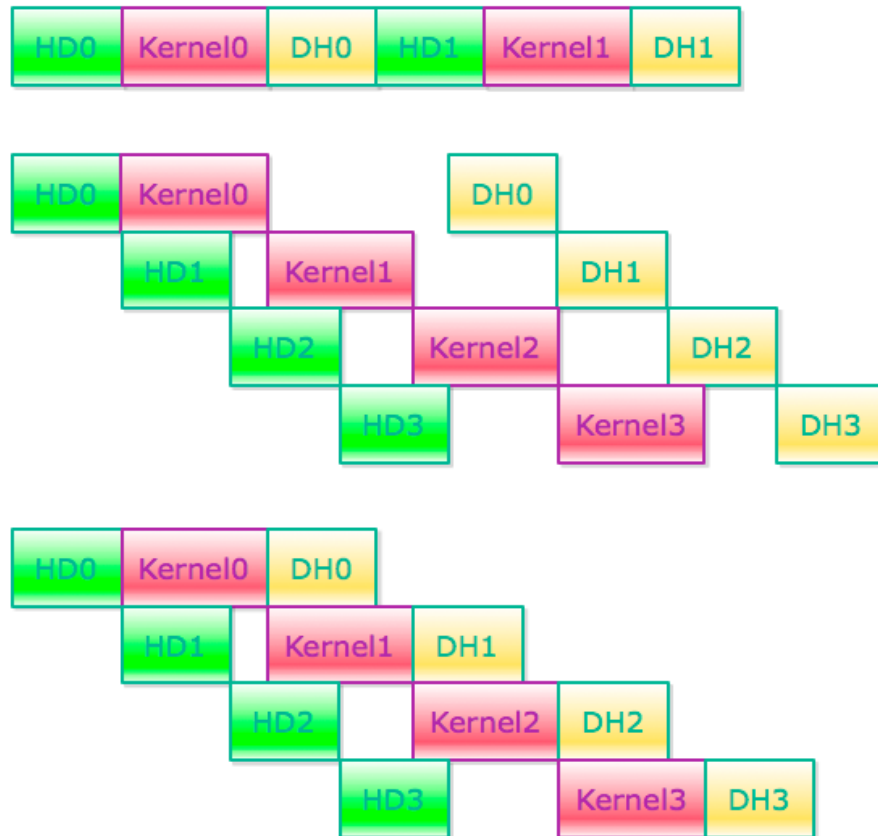


Fig. 3: Top: all operation in default stream. Mid: concurrent streams with one copy engine. Bottom: concurrent streams with two copy engine.