# Parallel Computing with GPUs

## Profiling

Dr Paul Richmond

http://paulrichmond.shef.ac.uk/teaching/COM4521/
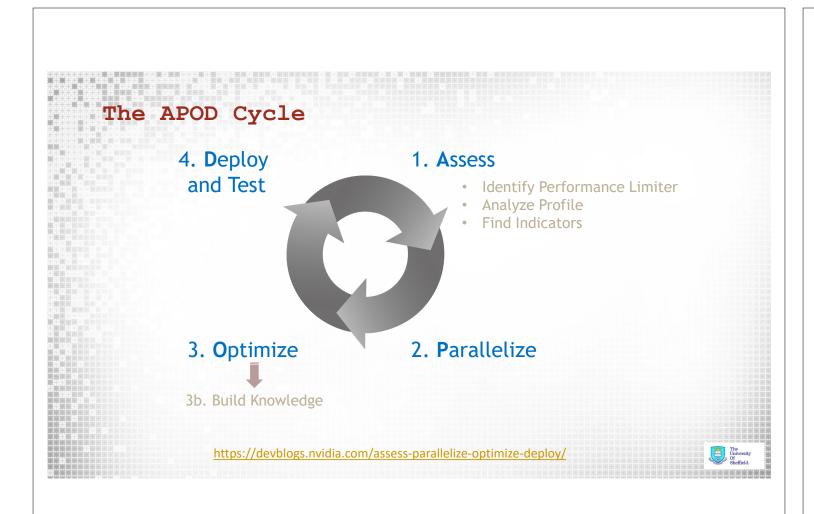
---

## Credits

❑ The code and much of the content from this lecture is based on the GTC2016 Talk by C. Angerer and J. Progsch (NVIDIA)

  ❑ S6112 – CUDA Optimisation with NVIDIA Nsight for Visual Studio

  ❑ Provided by NVIDIA with thanks to Joe Bungo

❑ Content has been adapted to use Visual Profiler Guided Analysis where possible

❑ Additional steps and analysis have been added

---

## Learning Objectives

❑ Understand the key performance metrics of GPU code.

❑ Understand profiling metrics and relate this to approaches which they have already learnt to address limiting factors in their code.

❑ Appreciate memory vs compute bound code and be able to recognise factors which contribute to this.

---
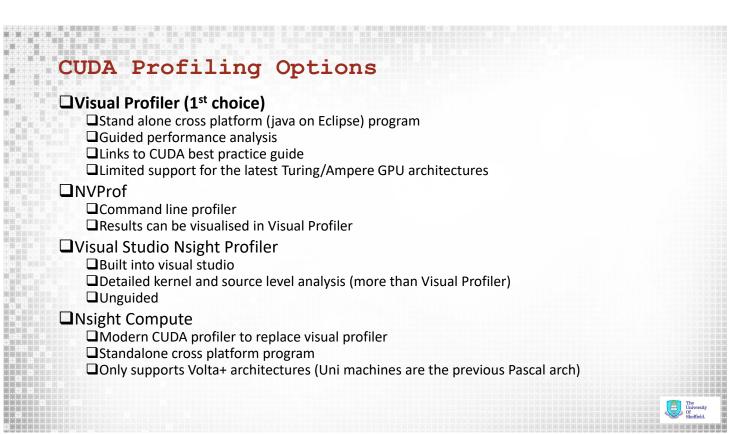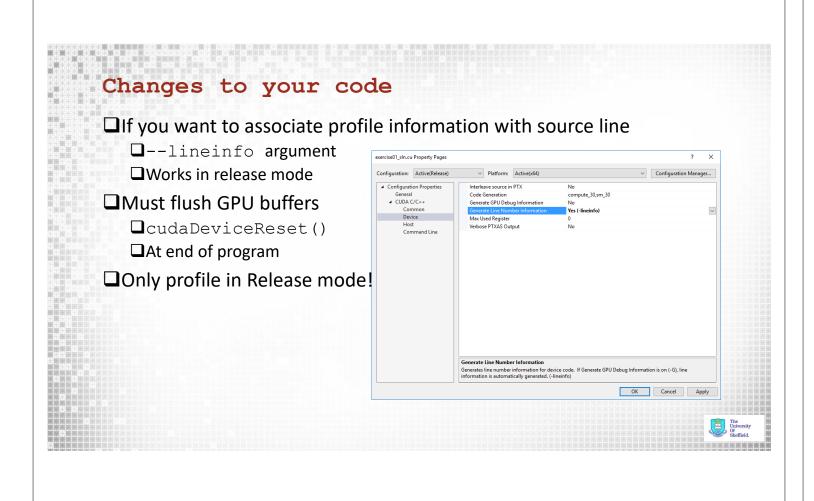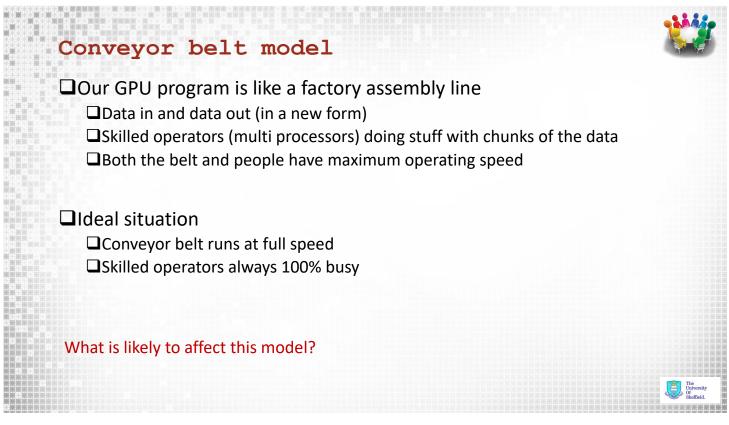
❑ Profiling Introduction

❑ The Problem

❑ Visual Profiler Guided Analysis

  ❑ Iteration 1

  ❑ Iteration 2

  ❑ Iteration 3

  ❑ Iteration 4
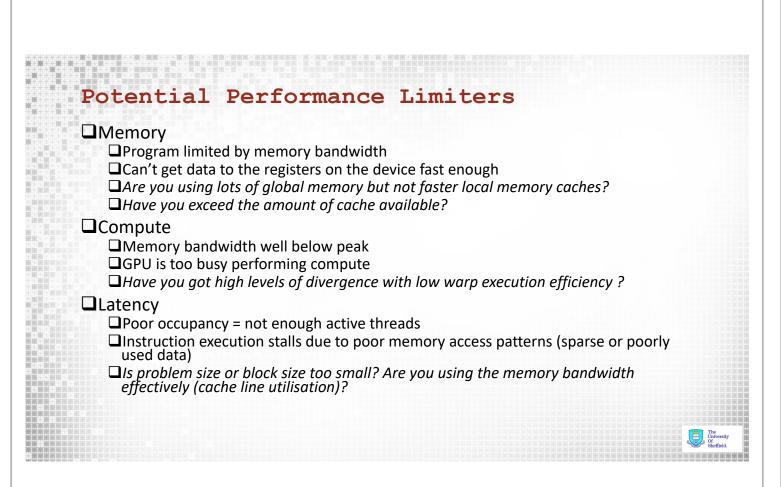
## The APOD Cycle

**4. Deploy and Test**

**1. Assess**
- Identify Performance Limiter
- Analyze Profile
- Find Indicators

**3. Optimize**

3b. Build Knowledge

**2. Parallelize**

https://devblogs.nvidia.com/assess-parallelize-optimize-deploy/

---

## CUDA Profiling Options

- ❑ **Visual Profiler (1st choice)**
  - ❑ Stand alone cross platform (java on Eclipse) program
  - ❑ Guided performance analysis
  - ❑ Links to CUDA best practice guide
  - ❑ Limited support for the latest Turing/Ampere GPU architectures
- ❑ NVProf
  - ❑ Command line profiler
  - ❑ Results can be visualised in Visual Profiler
- ❑ Visual Studio Nsight Profiler
  - ❑ Built into visual studio
  - ❑ Detailed kernel and source level analysis (more than Visual Profiler)
  - ❑ Unguided
- ❑ Nsight Compute
  - ❑ Modern CUDA profiler to replace visual profiler
  - ❑ Standalone cross platform program
  - ❑ Only supports Volta+ architectures (Uni machines are the previous Pascal arch)

---

## Changes to your code

- ❑ If you want to associate profile information with source line
  - ❑ `--lineinfo` argument
  - ❑ Works in release mode
- ❑ Must flush GPU buffers
  - ❑ `cudaDeviceReset()`
  - ❑ At end of program
- ❑ Only profile in Release mode!



---

## Conveyor belt model

- ❑ Our GPU program is like a factory assembly line
  - ❑ Data in and data out (in a new form)
  - ❑ Skilled operators (multi processors) doing stuff with chunks of the data
  - ❑ Both the belt and people have maximum operating speed

- ❑ Ideal situation
  - ❑ Conveyor belt runs at full speed
  - ❑ Skilled operators always 100% busy

What is likely to affect this model?

## Slide 1

# Potential Performance Limiters

- Memory
  - Program limited by memory bandwidth
  - Can't get data to the registers on the device fast enough
  - *Are you using lots of global memory but not faster local memory caches?*
  - *Have you exceed the amount of cache available?*
- Compute
  - Memory bandwidth well below peak
  - GPU is too busy performing compute
  - *Have you got high levels of divergence with low warp execution efficiency ?*
- Latency
  - Poor occupancy = not enough active threads
  - Instruction execution stalls due to poor memory access patterns (sparse or poorly used data)
  - *Is problem size or block size too small? Are you using the memory bandwidth effectively (cache line utilisation)?*

## Slide 2

- Profiling Introduction
- The Problem
- Visual Profiler Guided Analysis
  - Iteration 1
  - Iteration 2
  - Iteration 3
  - Iteration 4

## Slide 3
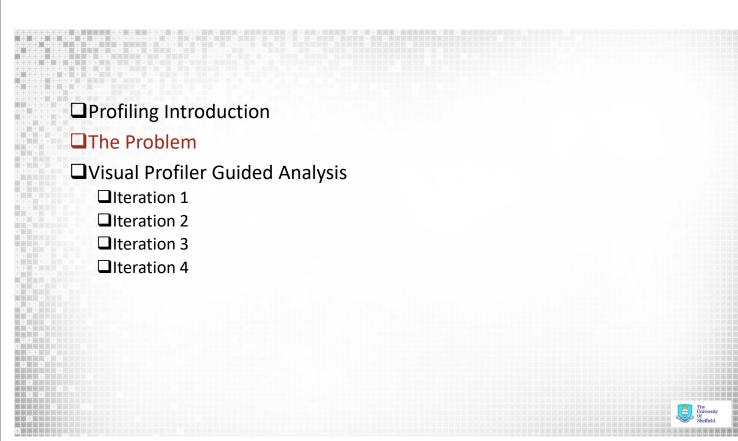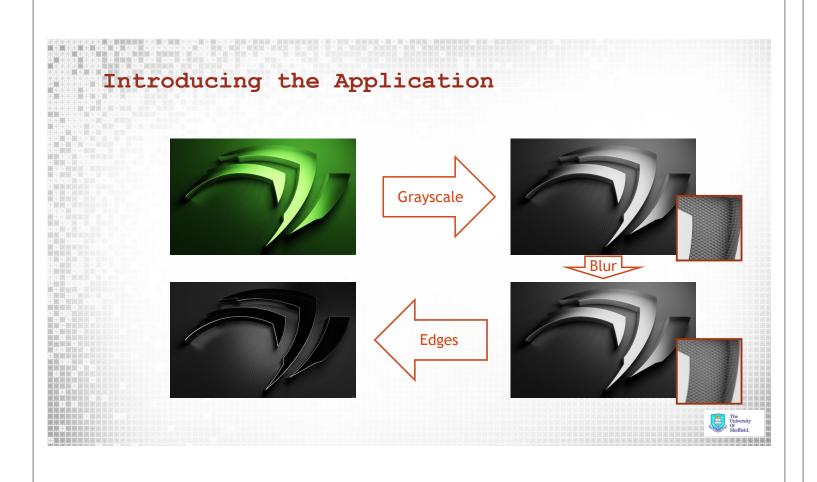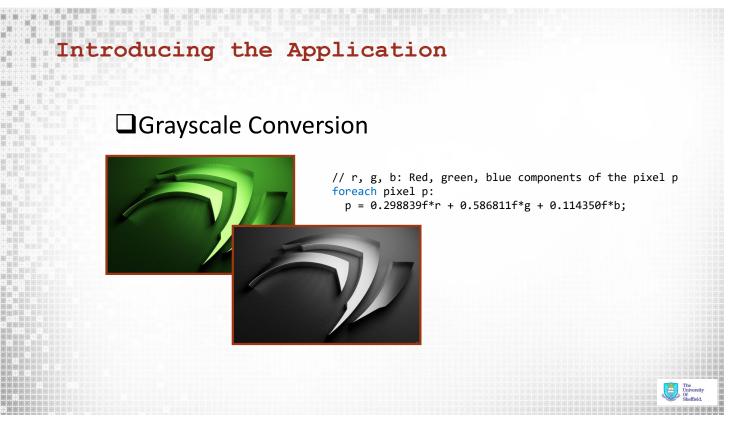
# Introducing the Application



## Slide 4

# Introducing the Application

- Grayscale Conversion

```
// r, g, b: Red, green, blue components of the pixel p
foreach pixel p:
    p = 0.298839f*r + 0.586811f*g + 0.114350f*b;
```

## Introducing the Application

### ❑ Blur: 7x7 Gaussian Filter



```
foreach pixel p:
  G = weighted sum of p and its 48 neighbors
  p  = G/256
```

*Image from Wikipedia*

---

## Introducing the Application

### ❑ Edges: 3x3 Sobel Filters



```
foreach pixel p:
  Gx = weighted sum of p and its 8 neighbors
  Gy = weighted sum of p and its 8 neighbors
  p  = sqrt(Gx + Gy)
```

Weights for Gx:

| -1 | 0 | 1 |
|----|---|---|
| -2 | 0 | 2 |
| -1 | 0 | 1 |

Weights for Gy:

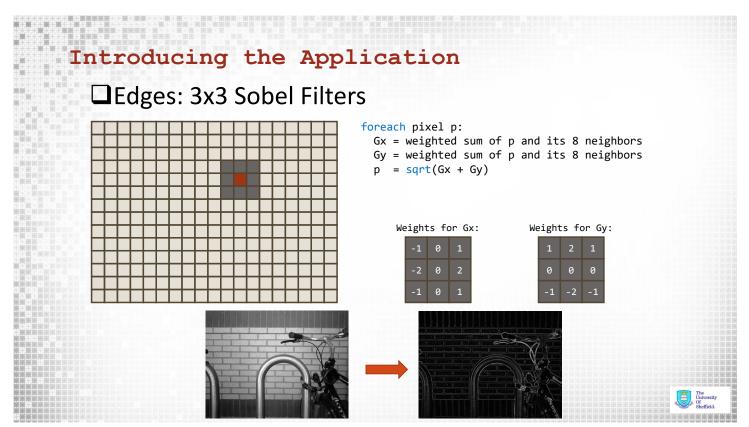| 1 | 2 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -2 | -1 |

---

## The Starting Code

```c
void gaussian_filter_7x7_v0(int w, int h, const uchar *src, uchar *dst)
{
  // Position of the thread in the image.
  const int x = blockIdx.x*blockDim.x + threadIdx.x;
  const int y = blockIdx.y*blockDim.y + threadIdx.y;

  // Early exit if the thread is not in the image.
  if( !in_img(x, y, w, h) )
    return;

  // Load the 48 neighbours and myself.
  int n[7][7];
  for( int j = -3 ; j <= 3 ; ++j )
    for( int i = -3 ; i <= 3 ; ++i )
      n[j+3][i+3] = in_img(x+i, y+j, w, h) ? (int) src[(y+j)*w + (x+i)] : 0;

  // Compute the convolution.
  int p = 0;
  for( int j = 0 ; j < 7 ; ++j )
    for( int i = 0 ; i < 7 ; ++i )
      p += gaussian_filter[j][i] * n[j][i];

  // Store the result.
  dst[y*w + x] = (uchar) (p / 256);
}
```

What is good and
what is bad?

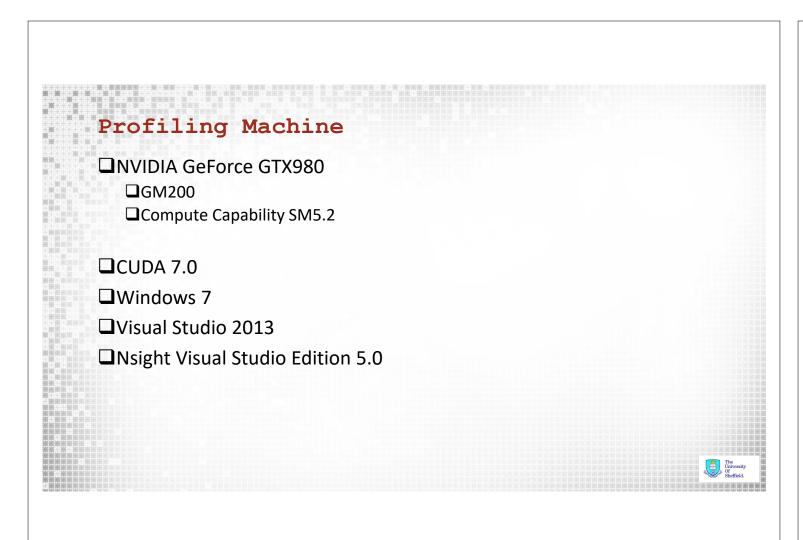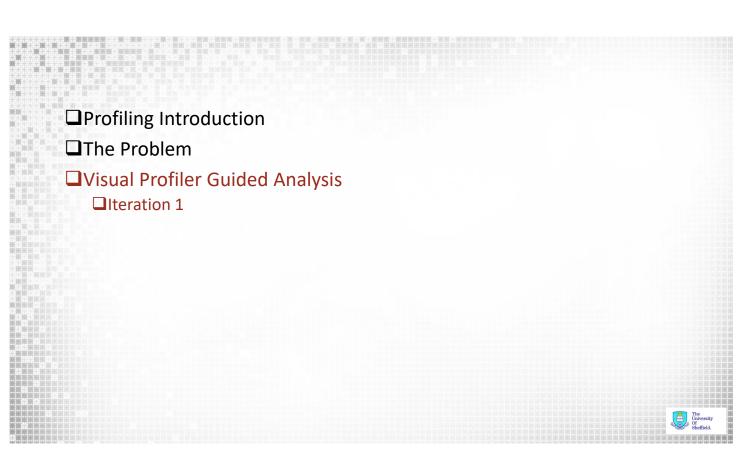❑ https://github.com/chmaruni/nsight-gtc

---

## The Starting Code

```c
void gaussian_filter_7x7_v0(int w, int h, const uchar *src, uchar *dst)
{
  // Position of the thread in the image.
  const int x = blockIdx.x*blockDim.x + threadIdx.x;
  const int y = blockIdx.y*blockDim.y + threadIdx.y;

  // Early exit if the thread is not in the image.
  if( !in_img(x, y, w, h) )
    return;

  // Load the 48 neighbours and myself.
  int n[7][7];
  for( int j = -3 ; j <= 3 ; ++j )
    for( int i = -3 ; i <= 3 ; ++i )
      n[j+3][i+3] = in_img(x+i, y+j, w, h) ? (int) src[(y+j)*w + (x+i)] : 0;

  // Compute the convolution.
  int p = 0;
  for( int j = 0 ; j < 7 ; ++j )
    for( int i = 0 ; i < 7 ; ++i )
      p += gaussian_filter[j][i] * n[j][i];

  // Store the result.
  dst[y*w + x] = (uchar) (p / 256);
}
```

What is good and
what is bad?

❑ https://github.com/chmaruni/nsight-gtc

# Profiling Machine

- NVIDIA GeForce GTX980
  - GM200
  - Compute Capability SM5.2

- CUDA 7.0
- Windows 7
- Visual Studio 2013
- Nsight Visual Studio Edition 5.0

---

- Profiling Introduction
- The Problem
- Visual Profiler Guided Analysis
  - Iteration 1

---



---



CUDA API Calls

Device Activity


Hints

# Results

⚠ **Low Kernel Concurrency** [ 0 ns / 5.94 ms = 0% ]

The percentage of time when two kernels are being executed in parallel is low.

More...

❑ We are using only a single stream

❑ Kernels have data dependencies so cant be executed in parallel

⚠ **Low Compute Utilization** [ 5.94 ms / 385.676 ms = 1.5% ]

The multiprocessors of one or more GPUs are mostly idle.

More...

❑ This is a problem

❑ The guided analysis will try and address this


Guided Analysis

guassian_filter_7x7_v0 kernel has highest rank



What is this telling us about our code?



Memory Bound Problem!



# Memory vs Compute vs Latency

60%

| Comp | Mem | Comp | Mem | Comp | Mem | Comp | Mem |

**Compute Bound** · **Bandwidth Bound** · **Latency Bound** · **Compute and Bandwidth Bound**

Better Occupancy might improve compute use

---

# What about occupancy?

❑Occupancy: *"number of active warps over max warps supported"*

❑Increasing achieved occupancy can hide latency

  ❑More warps available for execution = more to hide latency



■ The warp issues
□ The warp waits (latency)

Fully covered latency

Exposed latency, not enough warps

warp 0
warp 1
warp 2
warp 3
warp 4
warp 5
warp 6
warp 7
warp 8
warp 9

No warp issues

---

# Occupancy

❑In our case we are not achieving theoretical occupancy (we have latency)

What is the problem here?



■ The warp issues
□ The warp waits (latency)

warp 0
warp 1
warp 2
warp 3
warp 4
warp 5
warp 6
warp 7
warp 8
warp 9

---

# Occupancy

❑In our case we have good occupancy but still high latency

  ❑Schedulers cant find eligible warps at every cycle



■ The warp issues
□ The warp waits (latency)

Exposed latency at high occupancy

warp 0
warp 1
warp 2
warp 3
warp 4
warp 5
warp 6
warp 7
warp 8
warp 9

No warp issuing

Warps are waiting for memory (transactions)

**More information**



Transaction per access = 5:1
We are using only 20% of the effective bandwidth



# Transactions per access?

- Think back to Lecture 11
  - To get 100% efficiency **our threads need to access consecutive 4 byte values**
  - 32 Threads in warp accessing 4B each
    - 128B total via 4 L2 cache lines

```
__global__ void copy(float *odata, float* idata)
    int xid = blockIdx.x * blockDim.x + threadIdx.x;
    odata[xid] = idata[xid];
}
```

addresses from warp

Profiler is telling that we could use only 1 transaction but are using 4/5 (only 1 transaction required for each thread in warp to read a single byte char)



```
n[j+3][i+3] = in_img(x+i, y+j, w, h) ? (int) src[(y+j)*w + (x+i)] : 0;
```

Memory is indexed based on x==threadIdx.x: Suggests access is coalesced.
Cause not clear.....

## Analysis

❑ The limiting factor of our code is L2 Throughput
  ❑ There is nothing wrong with having high throughput
  ❑ Except: There is not enough compute to hide this
  ❑ We cant increase occupancy any further to hide this

❑ Solution: We need to reduce the time it takes to get data to the device to do compute on it. Either by
  ❑ Moving data closer to the SMPs
  ❑ **Making our L2 reads/writes more efficient**
    ❑ Currently ~4-5 Transactions/Access
    ❑ Our L2 cache lines are being used ineffectively

## Causes of Transaction per access: Striding?



```
__global__ void copy(float *odata, float* idata)
    int xid = (blockIdx.x * blockDim.x + threadIdx.x)* 2;
    odata[xid] = idata[xid];
}
```

❑ Lecture 11 example
  ❑ Strides (like above) cause poor transactions per access
  ❑ In the above case 8 transactions where we could have used 4

## Causes of Transaction per access: Offset?



| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 |

```
__global__ void copy(float *odata, float* idata)
    int xid = blockIdx.x * blockDim.x + threadIdx.x + 1;
    odata[xid] = idata[xid];
}
```

❑Lecture 11 Example:
  ❑If memory accesses are offset then parts of the cache line will be unused (shown in red) e.g.
  ❑Use thread blocks sizes of multiples of 32!

## What is our current data layout?

Blocks are 8x8



Block 1    Block 2    ...

Warps are 8x4

Warp 0

Warp 1

Why might this be a problem

## What is our current data layout?

Blocks are 8x8

Block 1    Block 2    ...

Warps are 8x4

Warp 0

threadIdx.x not consecutive within the warp

## Overfetch from L2 Cache

Line 245:              src[(y+j)*w + (x+i)]
Line 245 for i=0, j=0: src[x]              //threads 0-7 only



| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 |
| 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |
| 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 |
| 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 |

Warp 1
Warp 2

Data              Overfetch
0        8        16        24        32

Cache line (always aligned by 32B boundaries)

## Overfetch with L1 Caching

Line 245:            `src[(y+j)*w + (x+i)]`
Line 245 for i=0, j=0: `src[y*w + x]`



Data

Overfetch

Any Ideas for improving this?

---

## Optimisation: Improved Memory layout



- Minimum block width should be 32 (each thread requires only 1 byte)
- Use Layout of 32x2

Cache line (always aligned by 32B boundaries)

---

## Deploy: Improved Memory layout

| Kernel | Time (ms) | Speedup | Rel. Speedup |
|---|---|---|---|
| Gaussian_filter (Step 0) | 5.49 | 1.00x | - |
| **Gaussian_filter (Step 1a)** | **1.00** | **5.49x** | **5.49x** |

---

## Break

- What do we expect the analysis to look like next?
- Any ideas for what else may be required?

## Half time summary

❑The guided profiler will help us optimise the right thing

❑Hotspot tells us the most appropriate place to optimise

❑Performance Limiter tells us what to focus on to improve

❑Code may be Memory, Compute or Latency Bound

❑Improvements so far
- ❑Changed the access pattern (by changing block size)
- ❑Reduced memory dependencies?

---

❑Profiling Introduction

❑The Problem

❑Visual Profiler Guided Analysis
- ❑Iteration 1
- ❑Iteration 2

---

## Identify the hotspot

❑Examine GPU Usage in Visual Profiler

❑Examine Individual Kernels
- ❑Gaussian filter kernel still the highest rank

**i  Kernel Optimization Priorities**
The following kernels are ordered by optimization importance based on execution time and achieved occupancy. Optimization of higher ranked kernels (those that appear first in the list) is more likely to improve performance compared to lower ranked kernels.

| Rank | Description |
|------|-------------|
| 100 | [ 1 kernel instances ] gaussian_filter_7x7_v0(int, int, unsigned char const *, unsigned char*) |
| 29 | [ 1 kernel instances ] sobel_filter_3x3_v0(int, int, unsigned char const *, unsigned char*) |
| 14 | [ 1 kernel instances ] rgba_to_grayscale_kernel_v0(int, int, uchar4 const *, unsigned char*) |

---

## Performance Limiter

**i  Kernel Performance Is Bound By Memory Bandwidth**
For device "GeForce GTX 980" the kernel's compute utilization is significantly lower than its memory utilization. These utilization levels indicate that the performance of the kernel is most likely being limited by the memory system. For this kernel the limiting factor in the memory system is the bandwidth of the texture instruction units within the multiprocessors.

## Tex Instruction Units?

- What are texture instruction units and why might our code be using them?



## Tex Instruction Units?

- What are texture instruction units and why might our code be using them?
  - Hint:

```
void gaussian_filter_7x7_v1(int w,
                            int h,
                            const uchar *src,
                            uchar *dst)
```



## Tex Instruction Units?

- What are texture instruction units and why might out code be using them?

```
void gaussian_filter_7x7_v1(int w,
                            int h,
                            const uchar *src,
                            uchar *dst)
```

- Compiler is reading `src` as read-only through Unified L1/Read-Only (texture cache)



## Guided Bandwidth Analysis

⚠ GPU Utilization Is Limited By Memory Instruction Execution

The kernel's performance is potentially limited by the texture instruction units within the multiprocessors. These units are responsible for executing the instructions that result in accesses to memory. The table below shows the memory bandwidth used by this kernel for the various types of memory on the device.

*Optimization: Examine the compute analysis results for this kernel to determine how to reduce utilization and improve efficiency of the texture instruction units.* More...

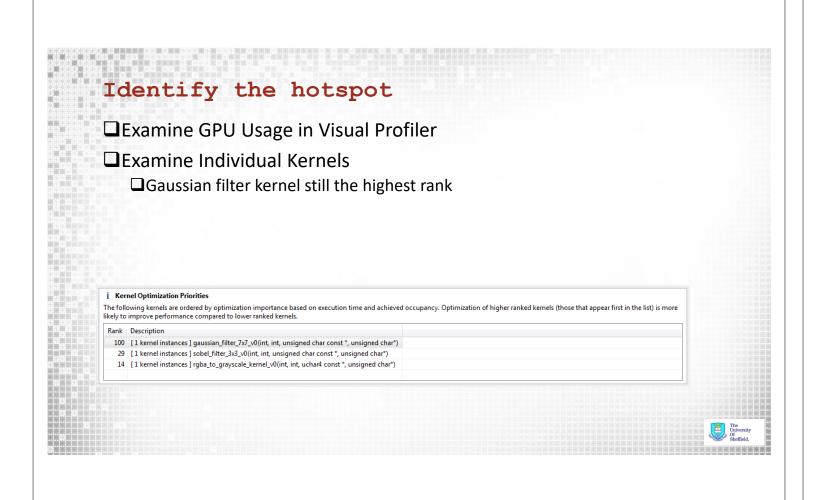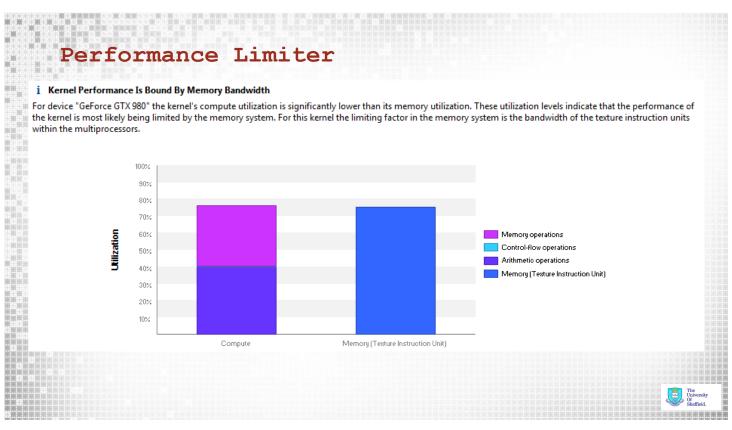| | Transactions | Bandwidth | Utilization |
|---|---|---|---|
| **Shared Memory** | | | |
| Shared Loads | 0 | 0 B/s | |
| Shared Stores | 0 | 0 B/s | |
| Shared Total | 0 | 0 B/s | |
| **L2 Cache** | | | |
| Reads | 11597012 | 415.583 GB/s | |
| Writes | 128006 | 4.587 GB/s | |
| Total | 11725018 | 420.17 GB/s | |
| **Unified Cache** | | | |
| Local Loads | 0 | 0 B/s | |
| Local Stores | 0 | 0 B/s | |
| Global Loads | 30364232 | 414.557 GB/s | |
| Global Stores | 128000 | 4.587 GB/s | |
| Texture Reads | 25061120 | 898.074 GB/s | |
| Unified Total | 55553352 | 1,317.218 GB/s | |
| **Device Memory** | | | |
| Reads | 157383 | 5.64 GB/s | |
| Writes | 127574 | 4.572 GB/s | |
| Total | 284957 | 10.212 GB/s | |
| **System Memory { PCIe configuration: Gen2 x16, 5 Gbit/s }** | | | |
| Reads | | 0 B/s | |
| Writes | 5 | 179.176 kB/s | |

- We are doing lots of reading/writing through unified cache

# Guided Bandwidth Analysis



☐ Still parts of the code reporting 2 transactions per access?

---

# Transaction per request

Line 245:               `src[(y+j)*w + (x+i)]`
Line 245 for **i=1**, j=0: `src[x+1]`

What is wrong with this access pattern?

---

# Transaction per request

Line 245:               `src[(y+j)*w + (x+i)]`
Line 245 for i=1, j=0: `src[x+1]`

What is wrong with this access pattern?

*Hint: Cache Lines are aligned by 32B boundaries*

---

# Transaction per request

Line 245:               `src[(y+j)*w + (x+i)]`
Line 245 for i=1, j=0: `src[x+1]`



We have an offset access pattern

# Slide 1: Guided Compute Analysis

**Guided Compute Analysis**

❑ The guided analysis suggests that lots of our compute cycles are spent issuing texture load/stores

ⓘ **Function Unit Utilization**

Different types of instructions are executed on different function units within each SM. Performance can be limited if a function unit is over-used by the instructions executed by the kernel. The following results show that the kernel's performance is not limited by overuse of any function unit.

Load/Store - Load and store instructions for shared and constant memory.
Texture - Load and store instructions for local, global, and texture memory.
Single - Single-precision integer and floating-point arithmetic instructions.
Double - Double-precision floating-point arithmetic instructions.
Special - Special arithmetic instructions such as sin, cos, popc, etc.
Control-Flow - Direct and indirect branches, jumps, and calls.



# Slide 2: Guided Latency Analysis: Occupancy

**Guided Latency Analysis: Occupancy**

⚠ **GPU Utilization May Be Limited By Register Usage**

Theoretical occupancy is less than 100% but is large enough that increasing occupancy may not improve performance. You can attempt the following optimization to increase the number of warps on each SM but it may not lead to increased performance.

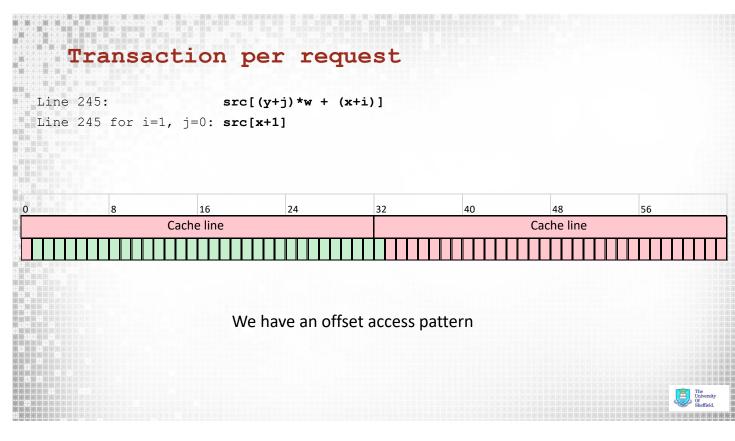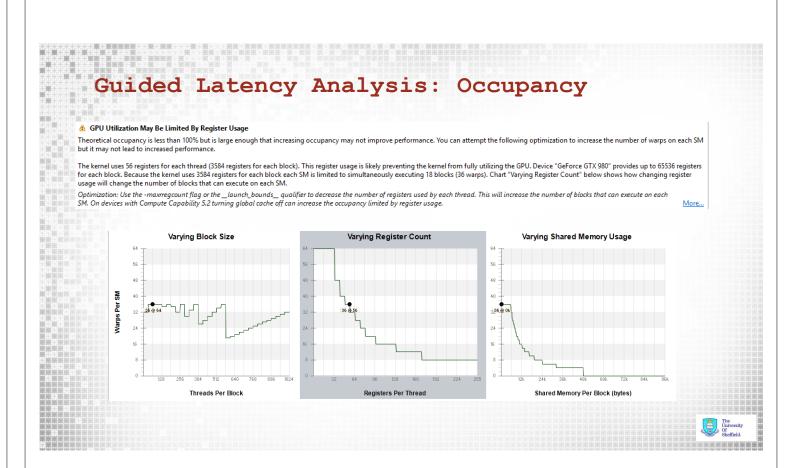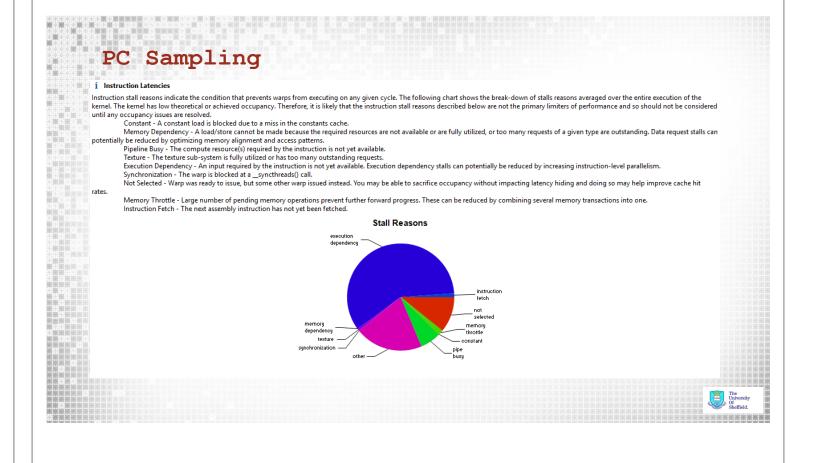The kernel uses 56 registers for each thread (3584 registers for each block). This register usage is likely preventing the kernel from fully utilizing the GPU. Device "GeForce GTX 980" provides up to 65536 registers for each block. Because the kernel uses 3584 registers for each block each SM is limited to simultaneously executing 18 blocks (36 warps). Chart "Varying Register Count" below shows how changing register usage will change the number of blocks that can execute on each SM.

*Optimization: Use the -maxrregcount flag or the __launch_bounds__ qualifier to decrease the number of registers used by each thread. This will increase the number of blocks that can execute on each SM. On devices with Compute Capability 5.2 turning global cache off can increase the occupancy limited by register usage.*

More...



# Slide 3: Guided Latency Analysis: Occupancy

❑ Register usage is very high

❑ Occupancy currently limited by register usage

❑ Increasing occupancy might not help us however as we are dominated by texture load stores
  - ❑ More work per SMP will just mean even more texture load stores!
  - ❑ We can confirm this by looking at the unguided analysis: Kernel Latency



# Slide 4: PC Sampling

**PC Sampling**

ⓘ **Instruction Latencies**

Instruction stall reasons indicate the condition that prevents warps from executing on any given cycle. The following chart shows the break-down of stalls reasons averaged over the entire execution of the kernel. The kernel has low theoretical or achieved occupancy. Therefore, it is likely that the instruction stall reasons described below are not the primary limiters of performance and so should not be considered until any occupancy issues are resolved.

Constant - A constant load is blocked due to a miss in the constants cache.
Memory Dependency - A load/store cannot be made because the required resources are not available or are fully utilized, or too many requests of a given type are outstanding. Data request stalls can potentially be reduced by optimizing memory alignment and access patterns.
Pipeline Busy - The compute resource(s) required by the instruction is not yet available.
Texture - The texture sub-system is fully utilized or has too many outstanding requests.
Execution Dependency - An input required by the instruction is not yet available. Execution dependency stalls can potentially be reduced by increasing instruction-level parallelism.
Synchronization - The warp is blocked at a __syncthreads() call.
Not Selected - Warp was ready to issue, but some other warp issued instead. You may be able to sacrifice occupancy without impacting latency hiding and doing so may help improve cache hit rates.
Memory Throttle - Large number of pending memory operations prevent further forward progress. These can be reduced by combining several memory transactions into one.
Instruction Fetch - The next assembly instruction has not yet been fetched.

## Execution/Memory Dependency

❑Rank these are best to worst

❑Which have instruction and memory dependencies?

```
int a = b + c;
int d = a + e;
//b, c and e are local ints
```

```
int a = b[i];
int d = a + e;
//b is global memory
//I and e are local ints
```

```
int a = b + c;
int d = e + f;
//b, c, e and f are local ints
```

## Instruction/Memory Dependency

❑Rank these are best to worst

❑Which have instruction and memory dependencies?

```
int a = b + c;
int d = a + e;
//b, c and e are local ints
```

```
int a = b[i];
int d = a + e;
//b is global memory
//i and e are local ints
```

```
int a = b + c;
int d = e + f;
//b, c, e and f are local ints
```

- ❑ Instruction Dependency
- ❑ Second add must wait for first

- ❑ Memory Dependency
- ❑ Second add must wait for memory request
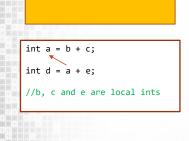
- ❑ No dependencies
- ❑ Independent Adds

## Analysis

❑Our compute engine is dominated by load/store instructions for the texture cache
  ❑Our texture bandwidth is good BUT
❑Our warps are stalling as instructions are waiting to issue texture fetch instructions
❑We still have poorly aligned access pattern within our inner loops

❑Solution: Reduce dependencies on texture loads
  ❑Move data closer to the SMP
  ❑Only read from global memory with nicely aligned cache lines
  ❑**How?**

## Analysis

❑Our compute engine is dominated by load/store instructions for the texture cache
  ❑Our texture bandwidth is good BUT
❑Our warps are stalling as instructions are waiting to issue texture fetch instructions
❑We still have poorly aligned access pattern within our inner loops

❑Solution: Reduce dependencies on texture loads
  ❑Move data closer to the SMP
  ❑Only read from global memory with nicely aligned cache lines
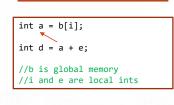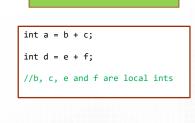  ❑**Shared Memory**

## Shared Memory

Single thread uses 7x7= 42 values

Use shared memory to store all pixels for the block

What important factor should we be considering?

Single block (32x4) uses 38x10 = 680 values

Also increased Block size

---

## Shared Memory

Single thread uses 7x7= 42 values

Use shared memory to store all pixels for the block

`__shared__ unsigned char smem_pixels[10][64]`

SM bank conflicts

Single block (32x4) uses 38x10 = 680 values

Also increased Block size

---

## BUT WAIT!!!!!!!!!!!!!!!!

❑ Wouldn't aligned char access have 4 way bank conflicts?
   ❑ NOT for Compute Mode 2.0+…

*"A shared memory request for a warp does not generate a bank conflict between two threads that access any address within the same 32-bit word (even though the two addresses fall in the same bank): In that case, for read accesses, the word is **broadcast** to the requesting threads (multiple words can be broadcast in a single transaction) …"*

| Thread | Bank |
| --- | --- |
| 0 | 0 |
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |
| 4 | 4 |
| 5 | 5 |
| 6 | 6 |
| 7 | 7 |
| | … |

I.e. A Stride of less than 1 (4B word) can be read conflict free if threads access aligned data

---

## Improvement

❑ Significant

| Kernel | Time (ms) | Speedup | Rel. Speedup |
| --- | --- | --- | --- |
| Gaussian_filter (Step 0) | 5.49 | 1.00x | - |
| Gaussian_filter (Step 1a) | 1.00 | 5.49x | 5.49x |
| **Gaussian_filter (Step 40)** | **0.49** | **11.20x** | **2.04x** |

## Slide 1

## Slide 2

# Identify the hotspot

❑Examine GPU Usage in Visual Profiler

❑Examine Individual Kernels

  ❑Gaussian filter kernel still the highest rank

    ❑Getting much closer though

**i Kernel Optimization Priorities**

The following kernels are ordered by optimization importance based on execution time and achieved occupancy. Optimization of higher ranked kernels (those that appear first in the list) is more likely to improve performance compared to lower ranked kernels.

| Rank | Description |
|------|-------------|
| 100 | [ 1 kernel instances ] gaussian_filter_7x7_v2(int, int, unsigned char const *, unsigned char*) |
| 60 | [ 1 kernel instances ] sobel_filter_3x3_v0(int, int, unsigned char const *, unsigned char*) |
| 29 | [ 1 kernel instances ] rgba_to_grayscale_kernel_v0(int, int, uchar4 const *, unsigned char*) |

## Slide 3

# Performance Limiter

**i Kernel Performance Is Bound By Memory Bandwidth**

For device "GeForce GTX 980" the kernel's compute utilization is significantly lower than its memory utilization. These utilization levels indicate that the performance of the kernel is most likely being limited by the memory system. For this kernel the limiting factor in the memory system is the bandwidth of the Shared memory.



❑Actually very close to magical 60% of compute

❑Lets examine

  ❑1) The compute analysis

  ❑2) The latency analysis

## Slide 4

# Guided Bandwidth Analysis

**⚠ GPU Utilization Is Limited By Memory Bandwidth**

The following table shows the memory bandwidth used by this kernel for the various types of memory on the device. The table also shows the utilization of each memory type relative to the maximum throughput supported by the memory. The results show that the kernel's performance is potentially limited by the bandwidth available from one or more of the memories on the device.

*Optimization: Try the following optimizations for the memory with high bandwidth utilization:*
*Shared Memory - If possible use 64-bit accesses to shared memory and 8-byte bank mode to achieved 2x throughput.*
*L2 Cache - Align and block kernel data to maximize L2 cache efficiency.*
*Unified Cache - Reallocate texture data to shared or global memory. Resolve alignment and access pattern issues for global loads and stores.*
*Device Memory - Resolve alignment and access pattern issues for global loads and stores.*
*System Memory (via PCIe) - Make sure performance critical data is placed in device or shared memory.* More...

| | Transactions | Bandwidth | Utilization |
|---|---|---|---|
| **Shared Memory** | | | |
| Shared Loads | 6272000 | 1,826.774 GB/s | |
| Shared Stores | 638720 | 186.033 GB/s | |
| Shared Total | 6910720 | 2,012.807 GB/s | High |
| **L2 Cache** | | | |
| Reads | 1318182 | 95.983 GB/s | |
| Writes | 128006 | 9.321 GB/s | |
| Total | 1446188 | 105.304 GB/s | |
| **Unified Cache** | | | |
| Local Loads | 0 | 0 B/s | |
| Local Stores | 0 | 0 B/s | |
| Global Loads | 3171255 | 92.366 GB/s | |
| Global Stores | 128000 | 9.32 GB/s | |
| Texture Reads | 2540993 | 185.022 GB/s | |
| Unified Total | 5840248 | 286.707 GB/s | |
| **Device Memory** | | | |
| Reads | 128446 | 9.353 GB/s | |
| Writes | 128008 | 9.321 GB/s | |
| Total | 256454 | 18.674 GB/s | |
| **System Memory  [ PCIe configuration: Gen2 x16, 5 Gbit/s ]** | | | |
| Reads | 0 | 0 B/s | |
| Writes | 5 | 364.073 kB/s | |

# Compute Analysis

- We are simply doing lots of compute
- Additional floating point operations graph shows no activity i.e. all of our instructions are Integer

What are all of these integer instructions?

- Export PDF Report
1. CUDA Application Analysis
2. Performance-Critical Kernels
3. Compute, Bandwidth, or Latency Bound
4. Compute Resources

GPU compute resources limit the performance of a kernel when those resources are insufficient or poorly utilized.

Show Kernel Profile - Instruction Execution

Rerun Analysis

---

# Compute Analysis by Line

- Selecting the CUDA function from compute analysis results allows a line by line breakdown
  - This will switch to unguided analysis

Also PTX instruction breakdown provided

---

# Guided Latency Analysis

Would changing the block size, register usage or amount of shared memory per block improve occupancy?

---

# Guided Latency Analysis

Line by Line Breakdown

## Latency Overview: Other 32.25%

❑Stall reason other generally means that there is no obvious action to improve performance

❑Other stall reasons may indicate either;
1. Execution unit is busy
   ❑ Solution: Potentially reduce use of low throughput integer operations if possible
2. Register bank conflicts : a compiler issue that can sometimes be made worst by heavy use of vector data types
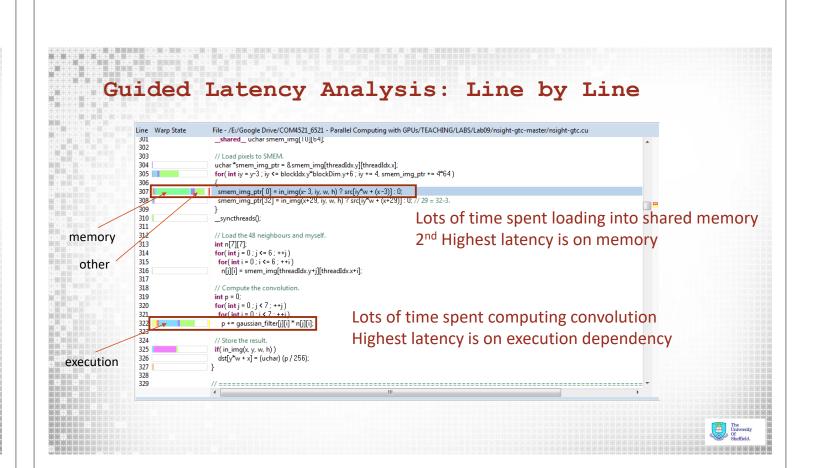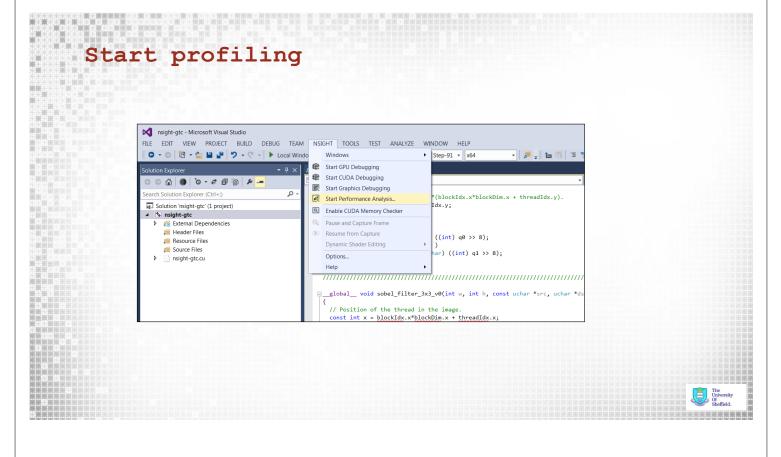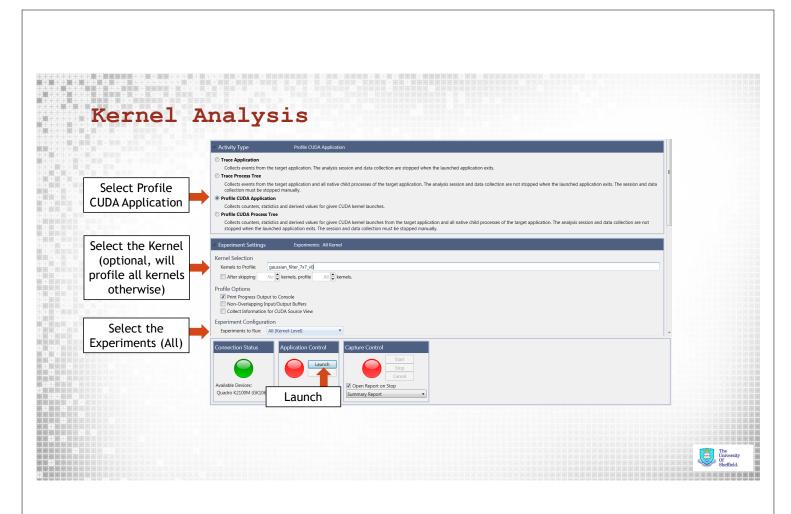   ❑ Solution: None
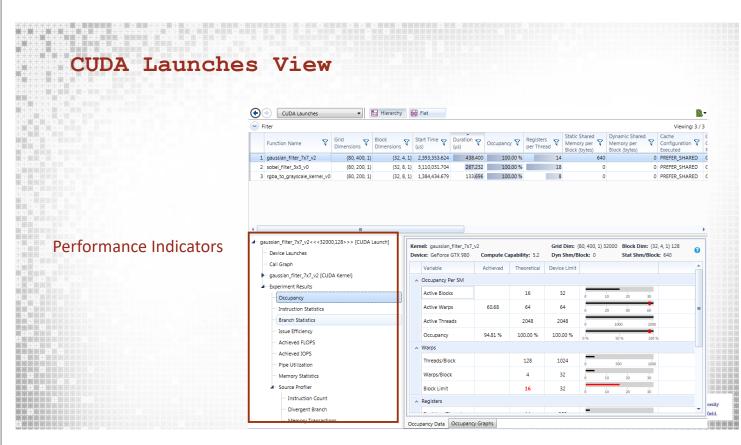3. Too few warps per scheduler
   ❑ Solution: Increase occupancy, decrease latency

---

## Guided Latency Analysis: Line by Line



memory

other

execution

Lots of time spent loading into shared memory
2nd Highest latency is on memory

Lots of time spent computing convolution
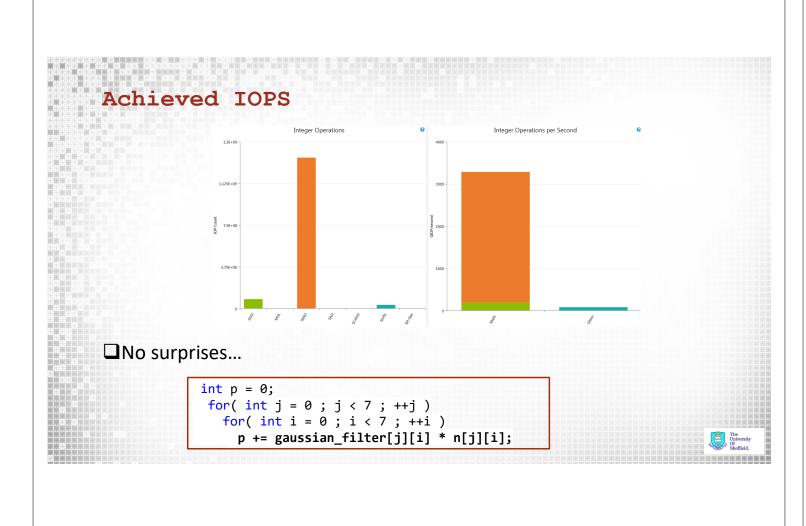Highest latency is on execution dependency

---

## 1st Analysis

❑We have a reasonably well balanced use of the from Compute and Memory pipes.

❑There is some latency in loading data to and from shared memory

❑Our compute cycles are dominated by Integer operations
   ❑What operations are they?
   ❑We can either examine the code and PTX instructions (from Compute or Latency Analysis) or run additional analysis via Nsight within Visual Studio
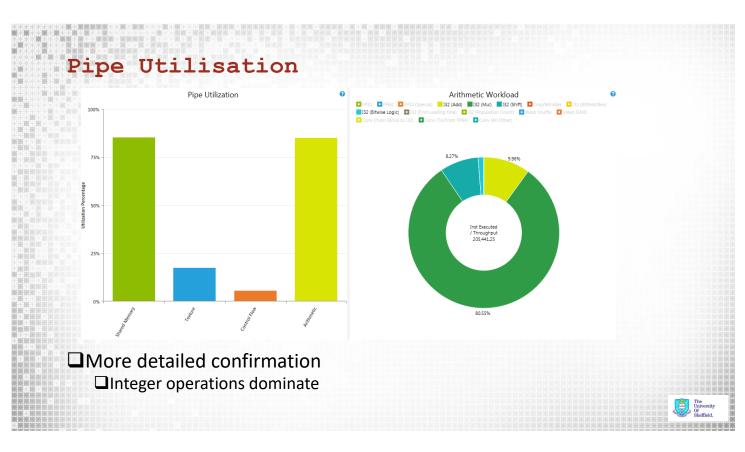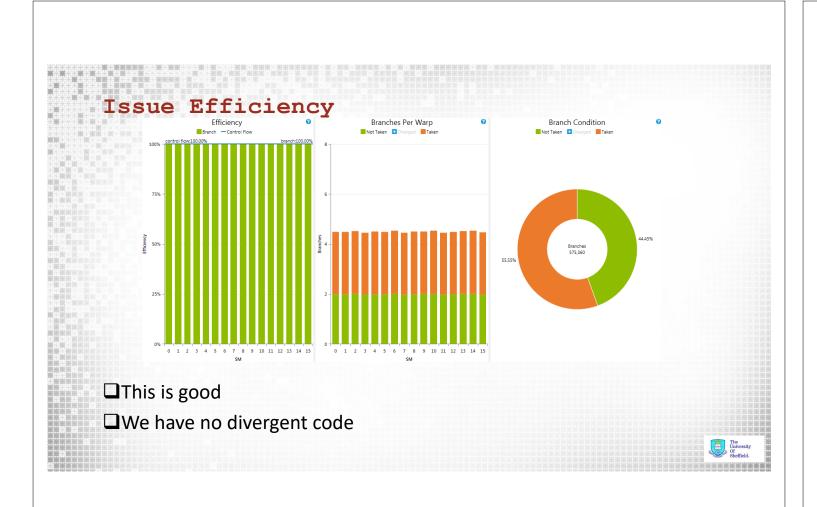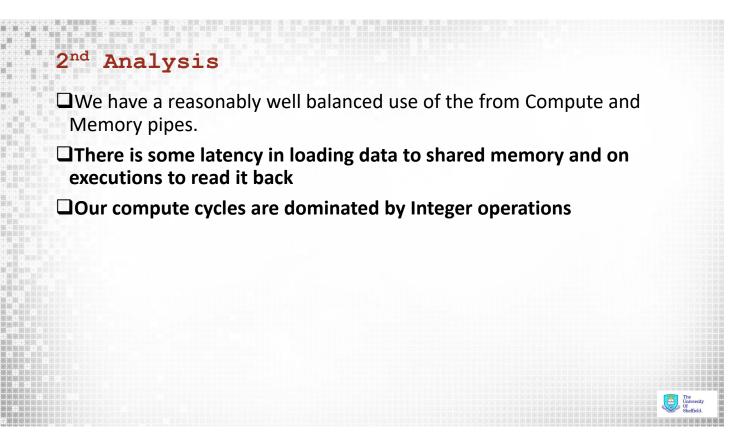      ❑More detailed analysis
      ❑Not guided like the visual profiler

---

## Start profiling

## Kernel Analysis



Select Profile CUDA Application

Select the Kernel (optional, will profile all kernels otherwise)

Select the Experiments (All)

Launch

## CUDA Launches View



Performance Indicators

## Achieved IOPS



☐ No surprises…

```
int p = 0;
 for( int j = 0 ; j < 7 ; ++j )
  for( int i = 0 ; i < 7 ; ++i )
   p += gaussian_filter[j][i] * n[j][i];
```

## Pipe Utilisation



☐ More detailed confirmation
  ☐ Integer operations dominate

## Issue Efficiency



- This is good
- We have no divergent code

---

## 2ⁿᵈ Analysis

- We have a reasonably well balanced use of the from Compute and Memory pipes.
- **There is some latency in loading data to shared memory and on executions to read it back**
- **Our compute cycles are dominated by Integer operations**

---

## There is some latency in loading data to shared memory and on executions to read it back

- Consider a simplified problem
- Each thread needs to load an r, g, b, a value into shared memory
  - **Which has fewer shared memory load instructions?**

```
__shared__ char sm[TPB*4];

char r,g,b,a;

r = sm[threadidx.x];
g = sm[threadidx.x+1];
b = sm[threadidx.x+2];
a = sm[threadidx.x+3];
```

```
__shared__ char4 sm[TPB];

char r,g,b,a;
char4 rgba;
rgba = sm[threadidx.x];
r = rgba.r;
g = rgba.g;
b = rgba.b;
a = rgba.a;
```

---

## There is some latency in loading data to shared memory and on executions to read it back

- Consider a simplified problem
- Each thread needs to load an r, g, b, a value into shared memory
  - **Which has fewer shared memory load instructions?**

```
__shared__ char sm[TPB*4];

char r,g,b,a;

r = sm[threadidx.x];
g = sm[threadidx.x+1];
b = sm[threadidx.x+2];
a = sm[threadidx.x+3];
```

```
__shared__ char4 sm[TPB];

char r,g,b,a;
char4 rgba;
rgba = sm[threadidx.x];
r = rgba.r;
g = rgba.g;
b = rgba.b;
a = rgba.a;
```

## Our compute cycles are dominated by Integer operations

```
int p = 0;
  for( int j = 0 ; j < 7 ; ++j )
    for( int i = 0 ; i < 7 ; ++i )
      p += gaussian_filter[j][i] * n[j][i];
```

❑Which of the following is faster?

```
int a, b, c;
a = sm_a[i]; b = sm_b[i];

c += a * b;
```

```
float a, b, c;
a = sm_a[i]; b = sm_b[i];

c += a * b;
```

---

Integer multiply add is 16 cycles

**Float combined multiply add is 4 cycles**

---

## Analysis

❑We have a reasonably well balanced use of the from Compute and Memory pipes.

❑There is some latency in loading data to shared memory and on executions to read it back
  ❑**Solution 1**: Reduce SM Load Stores dependencies by using wider requests. i.e. 4B values rather than 1B (chars)
  ❑I.e. Store shared memory values as 4B minimum

❑Our compute cycles are dominated by Integer operations
  ❑Almost all MAD operations
  ❑**Solution**: Change slower Integer MAD instructions to faster floating point FMAD instructions
  ❑I.e. Use floating point multiply and cast result to `uchar` at end

---

## Improvement

❑Significant

| Kernel | Time (ms) | Speedup | Rel. Speedup |
|---|---|---|---|
| Gaussian_filter (Step 0) | 5.49 | 1.00x | - |
| Gaussian_filter (Step 1a) | 1.00 | 5.49x | 5.49x |
| Gaussian_filter (Step 40) | 0.49 | 11.20x | 2.04x |
| **Gaussian_filter (Step 5a)** | **0.28** | **19.60x** | **1.75x** |

---

## Identify the hotspot

❑Examine GPU Usage in Visual Profiler

❑What should be our next step?

**i  Kernel Optimization Priorities**

The following kernels are ordered by optimization importance based on execution time and achieved occupancy. Optimization of higher ranked kernels (those that appear first in the list) is more likely to improve performance compared to lower ranked kernels.

| Rank | Description |
|------|-------------|
| 100 | [ 1 kernel instances ] sobel_filter_3x3_v0(int, int, unsigned char const *, unsigned char*) |
| 93 | [ 1 kernel instances ] gaussian_filter_7x7_v3_bis(int, int, unsigned char const *, unsigned char*) |
| 49 | [ 1 kernel instances ] rgba_to_grayscale_kernel_v0(int, int, uchar4 const *, unsigned char*) |

---

## Identify the hotspot

❑Examine GPU Usage in Visual Profiler

❑Examine Individual Kernels
    ❑Gaussian filter kernel no longer highest rank!
    ❑We can now optimise the sobel_filter kernel

**i  Kernel Optimization Priorities**

The following kernels are ordered by optimization importance based on execution time and achieved occupancy. Optimization of higher ranked kernels (those that appear first in the list) is more likely to improve performance compared to lower ranked kernels.

| Rank | Description |
|------|-------------|
| 100 | [ 1 kernel instances ] sobel_filter_3x3_v0(int, int, unsigned char const *, unsigned char*) |
| 93 | [ 1 kernel instances ] gaussian_filter_7x7_v3_bis(int, int, unsigned char const *, unsigned char*) |
| 49 | [ 1 kernel instances ] rgba_to_grayscale_kernel_v0(int, int, uchar4 const *, unsigned char*) |

❑Lets look at our Gaussian kernel anyway…

---

## Performance Limiter

**i  High Compute And Memory Utilization**

The kernel is utilizing greater than 80% of the available compute and memory performance of device "GeForce GTX 980". These utilization levels indicate that additional performance improvement may be difficult to achieved for the kernel.



❑Looking good

## VS NSight IOPS/ FLOPS Metrics



---

## Analysis

- ❑ Our algorithm is making good use of compute and memory
- ❑ Further improvement will be difficult (but not impossible)

- ❑ **Solution**: Optimise a different kernel
  - ❑ sobel_filter_kernel to get the same treatment
- ❑ **Solution**: Improve Gaussian kernel by changing the technique (parallelise differently)
  - ❑ Separable Filter: Compute horizontal and vertical convolution separately then approximate by binominal coefficients
  - ❑ Ensure we apply the same optimisations to separable filter version

---

## Improvement

| Kernel | Time (ms) | Speedup | Rel. Speedup |
|---|---|---|---|
| Gaussian_filter (Step 0) | 5.49 | 1.00x | - |
| Gaussian_filter (Step 1a) | 1.00 | 5.49x | 5.49x |
| Gaussian_filter (Step 40) | 0.49 | 11.20x | 2.04x |
| Gaussian_filter (Step 5a) | 0.28 | 19.60x | 1.75x |
| **Gaussian_filter (Step 9)** | **0.22** | **24.95x** | **1.27x** |

- ❑ 25x speedup on existing GPU code is pretty good

- ❑ Companion Code: https://github.com/chmaruni/nsight-gtc

---

## Summary

- ❑ Profiling with the Visual Profiler will give you guided analysis of how to improve your performance
  - ❑ Show you how to spot key metrics
- ❑ **We are trying to achieve good overall utilisation of the hardware (compute and memory engines)**
  - ❑ **Through an appreciation of memory and compute bounds**
- ❑ Follow the APOD cycle
  - ❑ Assess: What is the limiting factor, analyse and profile
  - ❑ Parallelise and improve (apply the knowledge you have learnt over the course)
  - ❑ Optimise
  - ❑ Deploy and Test
- ❑ If in doubt use the lab classes to seek guidance!