# Optimization Opportunities and Pitfalls when Implementing High Performance 2D Convolutions

# S4297

Ian Wainwright

High Performance Consulting Sweden

ian.wainwright@hpcsweden.se

# 2D Convolutions
# What are they?

### Filter

| 3 | 3 | 1 |
|---|---|---|
| 1 | 1 | 2 |
| 1 | 2 | 3 |

### Image

| 1 | 2 | 3 | 1 |
|---|---|---|---|
| 4 | 5 | 2 | 3 |
| 4 | 1 | 1 | 5 |
| 1 | 2 | 5 | 2 |

### Output

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

```
For all image pixels
    For all filter elements
        output += In(x,y)*Fil(x,y)
```

# 2D Convolutions
# What are they?

## Filter

| | | |
|---|---|---|
| 3 | 3 | 1 |
| 1 | 1 | 2 |
| 1 | 2 | 3 |

## Image

| | | | |
|---|---|---|---|
| 1 | 2 | 3 | 1 |
| 4 | 5 | 2 | 3 |
| 4 | 1 | 1 | 5 |
| 1 | 2 | 5 | 2 |

## Output

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |

```
For all image pixels
    For all filter elements
        output += In(x,y)*Fil(x,y)
```

# 2D Convolutions
# What are they?

### Filter

| | | |
|---|---|---|
| 3 | 3 | 1 |
| 1 | 1 | 2 |
| 1 | 2 | 3 |

### Image

| | | | |
|---|---|---|---|
| 1 * 3 | 2 * 3 | 3 * 1 | 1 |
| 4 * 1 | 5 * 1 | 2 * 2 | 3 |
| 4 * 1 | 1 * 2 | 1 * 3 | 5 |
| 1 | 2 | 5 | 2 |

### Output

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |

```
For all image pixels
    For all filter elements
        output += In(x,y)*Fil(x,y)
```

# 2D Convolutions
# What are they?

**Filter**

| | | |
|---|---|---|
| 3 | 3 | 1 |
| 1 | 1 | 2 |
| 1 | 2 | 3 |

**Image**

| | | | |
|---|---|---|---|
| 1 * 3 | 2 * 3 | 3 * 1 | 1 |
| 4 * 1 | 5 * 1 | 2 * 2 | 3 |
| 4 * 1 | 1 * 2 | 1 * 3 | 5 |
| 1 | 2 | 5 | 2 |

**Output**

| | | | |
|---|---|---|---|
| | | | |
| | 34 | | |
| | | | |
| | | | |

```
For all image pixels
    For all filter elements
        output += In(x,y)*Fil(x,y)
```

| | | |
|---|---|---|
| 3 | 6 | 3 |
| 4 | 5 | 4 |
| 4 | 2 | 3 |

∑ = 34

# 2D Convolutions

Tesla K20

GFLOPS = 3521 GFLOPS

BW 208 GB/s = 52 GigaFLOAT/s

→ 3521 GFLOPS/ 52 GFLOAT/s = *67* FLOPS / FLOAT is the theoretical break-even between bandwidth bound and compute bound.

The number of computations per output element in a 2D convolution is filter size * filter size * 2. → filter size = √(2 * *67*) = 11.5

→ break-even between compute and bandwidth bound.

*A 13\*13 2D convolution is in theory compute bound for the Tesla K20.*
Smaller sizes are bandwidth bound.

# The Tesla K20

Kepler 110 Whitepaper

**"192 single-precision floating point units"**

**"32 load/store units (LD/ST)"**

**→ 6 FMADs per LD/ST**

# The 2D Convolutions Implemented

**Filter**

- Not share filter
- Share filter

**Input image**

Filter directly from global

Filter via constant memory

Share filter via `__shfl`

Input via global memory

Store values per thread in regs

With and without `const __restrict__`

Single output per thread

Multiple output per thread

# The 2D Convolutions Implemented hpc

| Filter | | Input image | |

| No const __restrict__ | const __restrict__ on filter only | const __restrict__ on input only | const __restrict__ on both |

| No restrict | Filter restrict | Input restrict | Both restrict |

With and without const __restrict__

# The 2D Convolutions Implemented

# Everything in global memory

For all image pixels ⟶ Map to each thread
    For all filter pixels ⟶ For loop for each thread
        output **+= In(x,y)\*Fil(x,y)**

# Everything in global memory

For all image pixels ——————→ Map to each thread
   For all filter pixels ————→ For loop for each thread
     output += In(x,y)*Fil(x,y)   **GK 110: 6 FMADs per LD/ST**

    **FMAD**   **LD**   **LD** ——————→ **2 LD per FMAD**

| FPU1 | FPU2 | FPU3 | FPU4 | FPU5 | FPU6 | LD/ST |
|------|------|------|------|------|------|-------|
|      |      |      |      |      |      | LD In |
|      |      |      |      |      |      | LD Fil |
| FMAD |      |      |      |      |      | LD In |
|      |      |      |      |      |      | LD Fil |
|      | FMAD |      |      |      |      |       |

**At best we will utilize 1/12 of the hardware**

# Everything in global memory

**Roughly 4% utilization**



Input restrict

Both restrict

Filter restrict

No restrict

# The 2D Convolutions Implemented

# The 2D Convolutions Implemented

# Filter in constant, image in global

Roughly 7% utilization



Input restrict

No restrict

# The 2D Convolutions Implemented

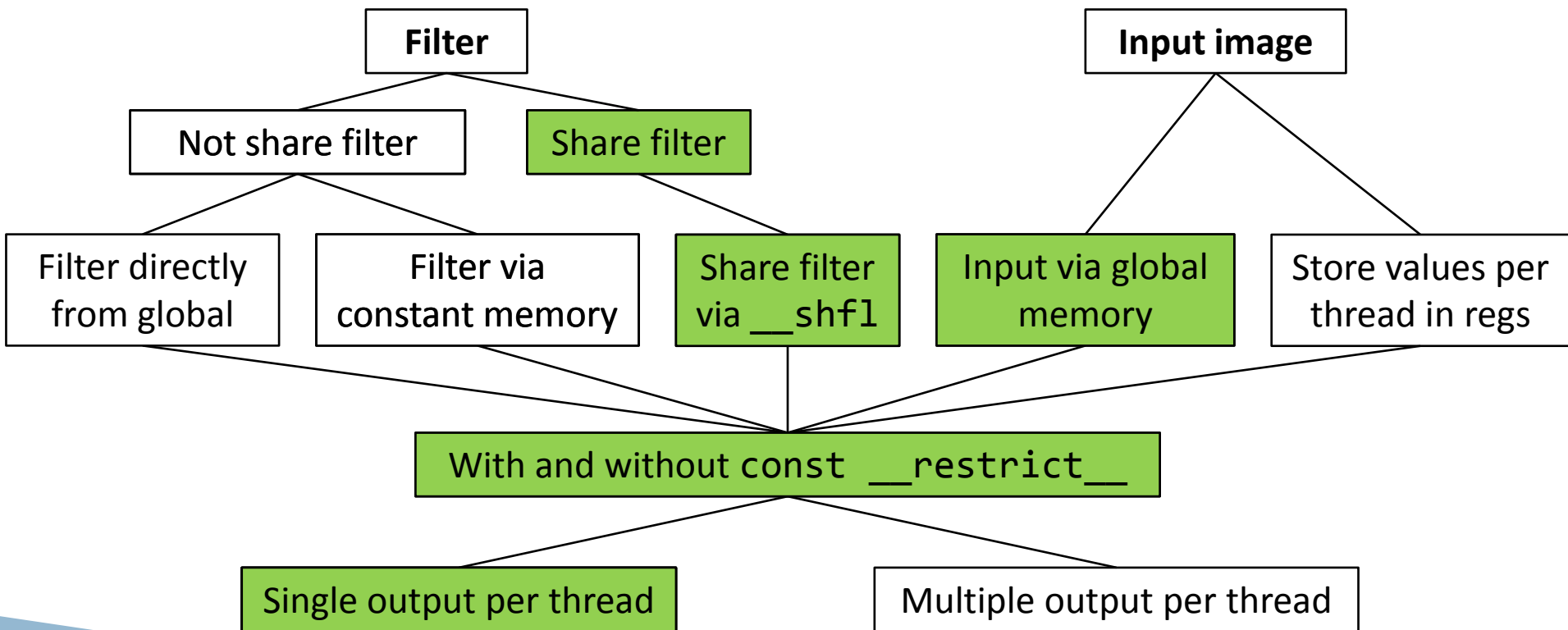# The 2D Convolutions Implemented
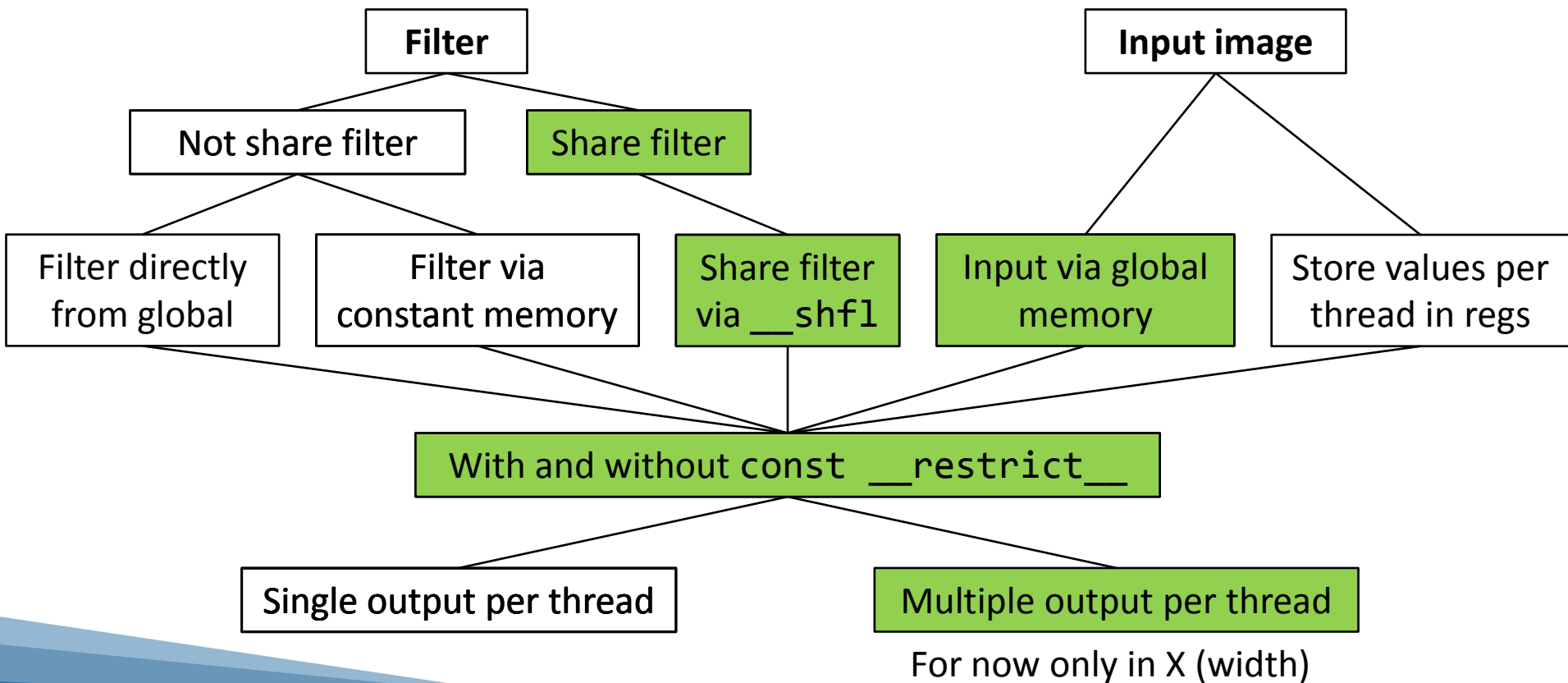
# Share filter via __shfl, image in global

# The 2D Convolutions Implemented

# The 2D Convolutions Implemented

**Filter**

**Input image**

Not share filter

Share filter

Filter directly from global

Filter via constant memory

Share filter via `__shfl`

Input via global memory

Store values per thread in regs

With and without `const __restrict__`

Single output per thread

Multiple output per thread

For now only in X (width)

# The 2D Convolutions Implemented



Filter

Not share filter — Share filter

Filter directly from global — Filter via constant memory — Share filter via __shfl

Input image

Input via global memory — Store values per thread in regs

With and without const __restrict__

Single output per thread — Multiple output per thread

For now only in X (width)

# Multiple outputs per thread

For all image pixels $\longrightarrow$ Map to each thread
    For all filter elements $\longrightarrow$ For loop for each thread
        output += **In(x.y)\*Fil(x,y)**


For all image pixels $\longrightarrow$ Map one thread to *several* outputs
    For all filter elements $\longrightarrow$ For loop for each thread
        *fil = Fil(x,y)*
        output0 += In(x+0,y)\**fil*
        output1 += In(x+1,y)\**fil*
        output2 += In(x+2,y)\**fil* ...

# Multiple outputs per thread

```
For all image pixels
    For all filter elements
        fil = Fil(x,y)
        output0 += In(x+0,y)*fil
        output1 += In(x+1,y)*fil
        output2 += In(x+2,y)*fil ...
```

**GK 110: 6 FMADs per LD/ST**

Place filter value in register, i.e. only 1 LD

**1 LD per FMAD**

| FPU1 | FPU2 | FPU3 | FPU4 | FPU5 | FPU6 | LD/ST |
|------|------|------|------|------|------|-------|
| | | | | | | LD Fil |
| | | | | | | LD Fil |
| | | | | | | LD In |
| FMAD | | | | | | LD In |
| | FMAD | | | | | LD In |
| FMAD | | | | | | LD In |
| | FMAD | | | | | |

**At best we will utilize 1/6 of the hardware**

# Filter in constant, image in global, multiple outputs per thread

*Still no reuse of input data*

# The 2D Convolutions Implemented

# The 2D Convolutions Implemented

**Filter**

Not share filter

Share filter

Filter directly from global

Filter via constant memory

Share filter via __shfl

**Input image**

Input via global memory

Store values per thread in regs

With and without const __restrict__

Single output per thread

Multiple output per thread

```
For all image pixels
    For all filter elements
        fil0 = Fil(x+0,y);
        in0 = In(x+0,y); in1 = In(x+1,y)...
        output0 += in0*fil0
        output1 += in1*fil0
        output2 += in2*fil0 ...
```

Load filter into register.
Load all input elements into register.
Do all FMAs.

| FPU1 | FPU2 | FPU3 | FPU4 | FPU5 | FPU6 | LD/ST |
|------|------|------|------|------|------|-------|
|      |      |      |      |      |      | LD filter |
|      |      |      |      |      |      | LD all In |
| FMAD |      |      |      |      |      | LD filter |
| FMAD |      |      |      |      |      | LD all In |
| FMAD | FMAD |      |      |      |      | LD filter |
| FMAD | FMAD |      |      |      |      | LD all In |
| FMAD | FMAD | FMAD |      |      |      |       |

**Filter-size*2 operations per load**

# Filter via `const`, reusing thread local input, multiple outputs per thread



No restrict

Input restrict

From 3 to 30 % utilization by storing input values in registers *if input image uses* `const __restrict__`

# The 2D Convolutions Implemented

**Filter**

- Not share filter
- Share filter

**Input image**

Filter directly from global

Filter via constant memory

Share filter via `__shfl`

Input via global memory

Store values per thread in regs

With and without `const __restrict__`

Single output per thread

Multiple output per thread

# The 2D Convolutions Implemented

**Filter**

**Input image**

Not share filter

Share filter

Filter directly from global

Filter via constant memory

Share filter via `__shfl`

Input via global memory

Store values per thread in regs

With and without `const __restrict__`

Single output per thread

Multiple output per thread

# Sharing filter via shfl, reusing thread local input, multiple outputs per thread

# Sharing filter via shfl, reusing thread local input, multiple outputs per thread



Input restrict

What does the code look like?

Both restrict

From 3 to 30 % utilization by storing input values in registers
*if input image uses*
`const __restrict__`

```
For all image pixels
    For all filter elements
        fil0 = Fil(x+0,y);
        in0 = In(x+0,y); in1 = In(x+1,y)...
        output0 += in0*fil0
        output1 += in1*fil0
        output2 += in2*fil0 ...
```

Load filter into register.
Load all input elements into register.
Do all FMAs.

| FPU1 | FPU2 | FPU3 | FPU4 | FPU5 | FPU6 | LD/ST |
|------|------|------|------|------|------|-------|
|      |      |      |      |      |      | LD filter |
|      |      |      |      |      |      | LD all In |
| FMAD |      |      |      |      |      | LD filter |
| FMAD |      |      |      |      |      | LD all In |
| FMAD | FMAD |      |      |      |      | LD filter |
| FMAD | FMAD |      |      |      |      | LD all In |
| FMAD | FMAD | FMAD |      |      |      |      |

**Filter-size*2 operations per load**

4 lines of code for reading input into register
9 lines of code for filter loop (shown below)
9 lines of code for output clean-up
= 22 lines of code. → *Maintainable*

```
121    for(i=0; i<T_FILTER_WIDTH; i++) {                                              1533312
122        if(threadIdx.x < T_FILTER_WIDTH) {                                          766656
123            the_warps_filter_coefs = filter_coefs[threadIdx.x + i*T_FILTER_WIDTH];  1916640
124        }
125
126        const int input_data_float_offset = line_width * y_input_pos + x_block_pos + (i - half_filter_width) * line_width…
127
128        load_data_float<nb_float_input_values>(my_input_values, input_data, input_data_float_offset);
129
130        for(j=0; j<T_FILTER_WIDTH; j++) {
131            // Shuffle to get the filter value from thread j;
132            present_filter_coef = __shfl(the_warps_filter_coefs, j);
133
134            for(int output_index=0; output_index<T_NB_OUTPUTS_PER_THREAD; output_index++) {
135                my_X_output[output_index] += present_filter_coef*my_input_values[output_index + j];   44466048 14
136            }
137        }
138    }
```

No hand-coded assembly
No explicit use of textures
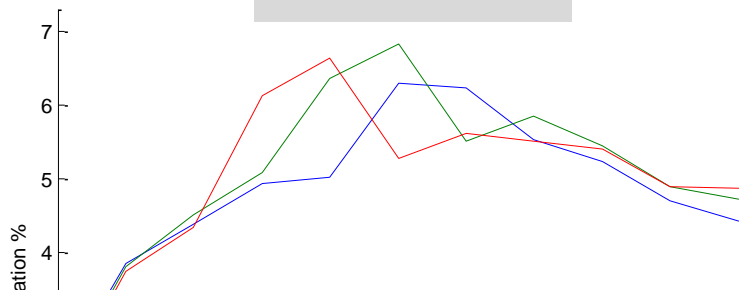No use of constant memory

Instead:
C++ Templates, const __restrict__, __shfl

**~ 100 FMAD**

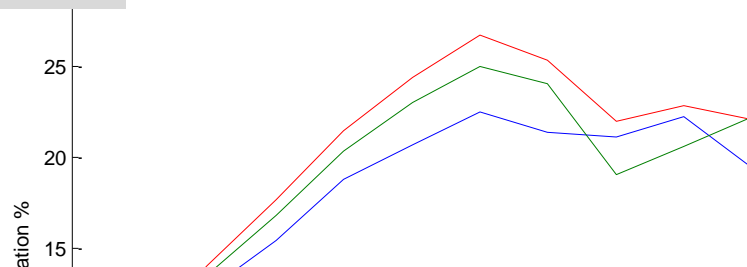# Use compile time in inner loops

Inner most loop is
*run time*

Both restrict

A factor of 3 difference in performance!

Inner most loop is
*compile time*

Both restrict



```
for(j=0; j<filtersize_compile_time_constant; j++)
{
    const float present_filter_coef = __shfl(the_warps_filter_coefs, j);

    for(int output_index=0; output_index<T_NB_OUTPUTS_PER_THREAD; output_index++)
    {
        my_X_output[output_index] += present_filter_coef*my_input_values[output_index + j];
    }
}
```
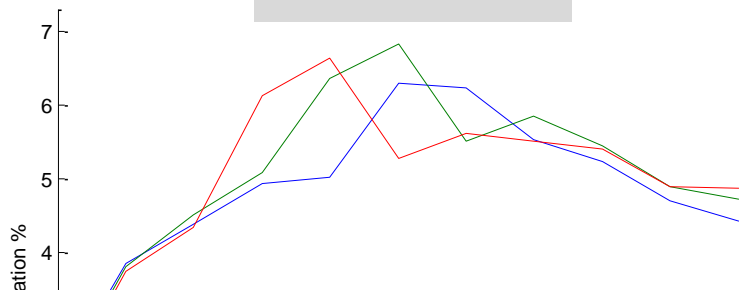
# Use compile time in inner loops
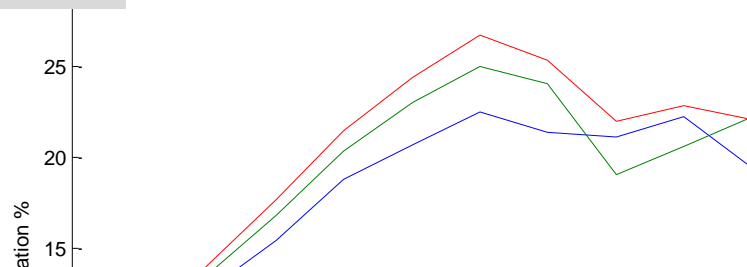
Inner most loop is *run time*

Both restrict

A factor of 3 difference in performance!

Inner most loop is *compile time*

Both restrict

```
for(j=0; j<filtersize_run_time_constant; j++)
{
    const float present_filter_coef = __shfl(the_warps_filter_coefs, j);

    for(int output_index=0; output_index<T_NB_OUTPUTS_PER_THREAD; output_index++)
    {
        my_X_output[output_index] += present_filter_coef*my_input_values[output_index + j];
    }
}
```

# 2D Convolutions
# What are they?

### Filter

| | | |
|---|---|---|
| 3 | 3 | 1 |
| 1 | 1 | 2 |
| 1 | 2 | 3 |

### Image

| | | | |
|---|---|---|---|
| 1 * 3 | 2 * 3 | 3 * 1 | 1 |
| 4 * 1 | 5 * 1 | 2 * 2 | 3 |
| 4 * 1 | 1 * 2 | 1 * 3 | 5 |
| 1 | 2 | 5 | 2 |

### Output

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |

# 2D Convolutions

hpc

### Filter

| | | |
|---|---|---|
| 3 | 3 | 1 |
| 1 | 1 | 2 |
| 1 | 2 | 3 |

### Image

| | | | |
|---|---|---|---|
| 1 * 3 | 2 * 3 | 3 * 1 | 1 |
| 4 * 1 | 5 * 1 | 2 * 2 | 3 |
| 4 * 1 | 1 * 2 | 1 * 3 | 5 |
| 1 | 2 | 5 | 2 |

### Output

| | | | |
|---|---|---|---|
| | | | |
| | 34 | | |
| | | | |
| | | | |

## Extend reuse of input data to Y also.

| | | |
|---|---|---|
| 3 | 6 | 3 |
| 4 | 5 | 4 |
| 4 | 2 | 3 |

∑ = 34

# Reuse input in X *and* Y

```
For all image pixels    ─────────→  Map one thread to several outputs
    For all filter pixels ─────────→  For loop for each thread
        fil = Fil(x,y)
        output0 += In(x+0,y)*fil
        output1 += In(x+1,y)*fil
        output2 += In(x+2,y)*fil
        output3 += In(x+3,y)*fil
        ...

For all image pixels    ─────────→  Map one thread to several outputs in x and y
    For all filter pixels ─────────→  For loop for each thread
        fil = Fil(x,y)
        out00 = In(x+0,y+0)*fil
        out01 = In(x+0,y+1)*fil
        out10 = In(x+1,y+0)*fil
        out11 = In(x+1,y+1)*fil
```
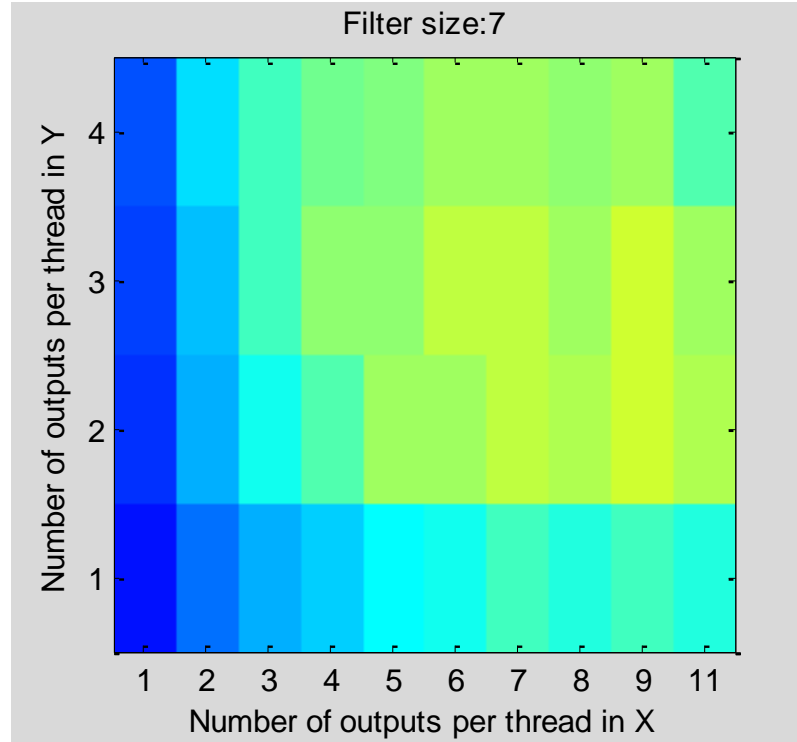
# Reuse input in X *and* Y
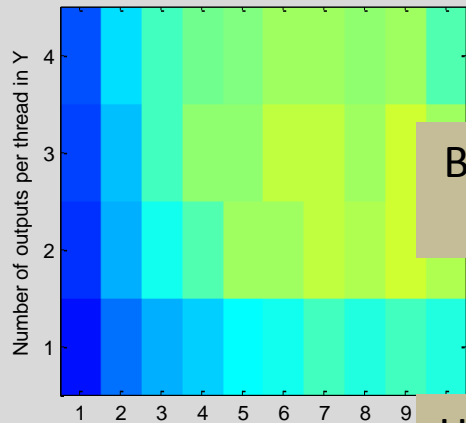


**Multiple outputs per thread in Y (height)**
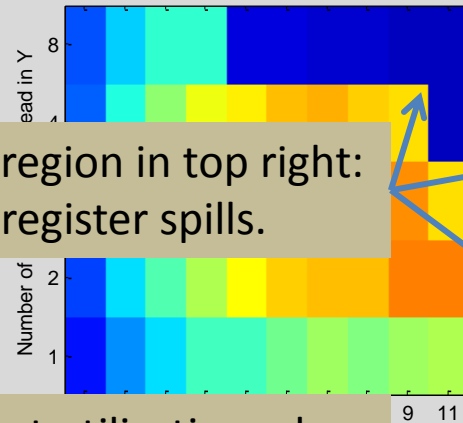
**Multiple outputs per thread in X (width)**

**Utilization heat map: blue is low utilization, red is high utilization.**
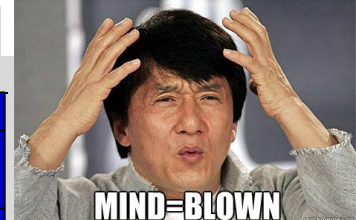
Filter size 7

Filter size 9

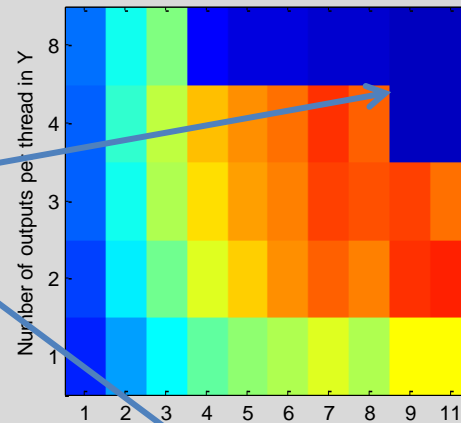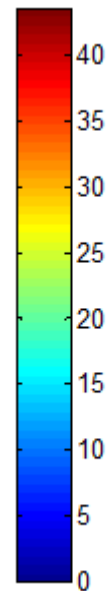Filter size 11

MIND=BLOWN

Utilization %

Blue region in top right: register spills.
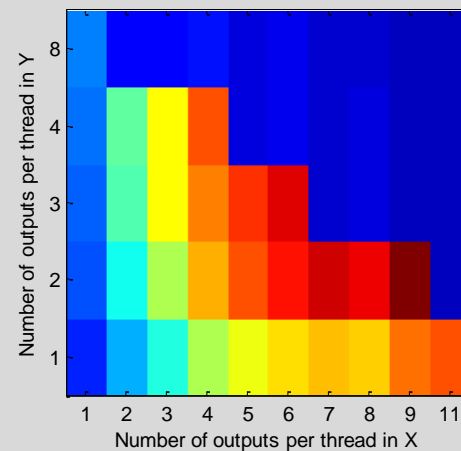
Highest utilization when we have 2 outputs in Y per thread

Filter size 13

Filter size 17

# 2D Convolutions Conclusion

DOs:
- Map multiple outputs to each thread.
- Use templates to hardcode loops as non-constant indexed arrays "[are] likely to [be] placed in local memory".
- Its helpfull to have a basic understanding and model of the hardware you're working with.
- Keep looking at you assembly: What lines map to register based compute, and what is LD/ST integer spaghetti code?

DON'Ts:

- Use run-time sizes in inner most loops.
- Use textures or constant memory.
const __restrict__ gets the job done and its very simple to use!

## Any questions?

High Performance Consulting Sweden

ian.wainwright@hpcsweden.se