

GPU optimization techniques and tools

A collection of ideas to maybe improve your GPU performance

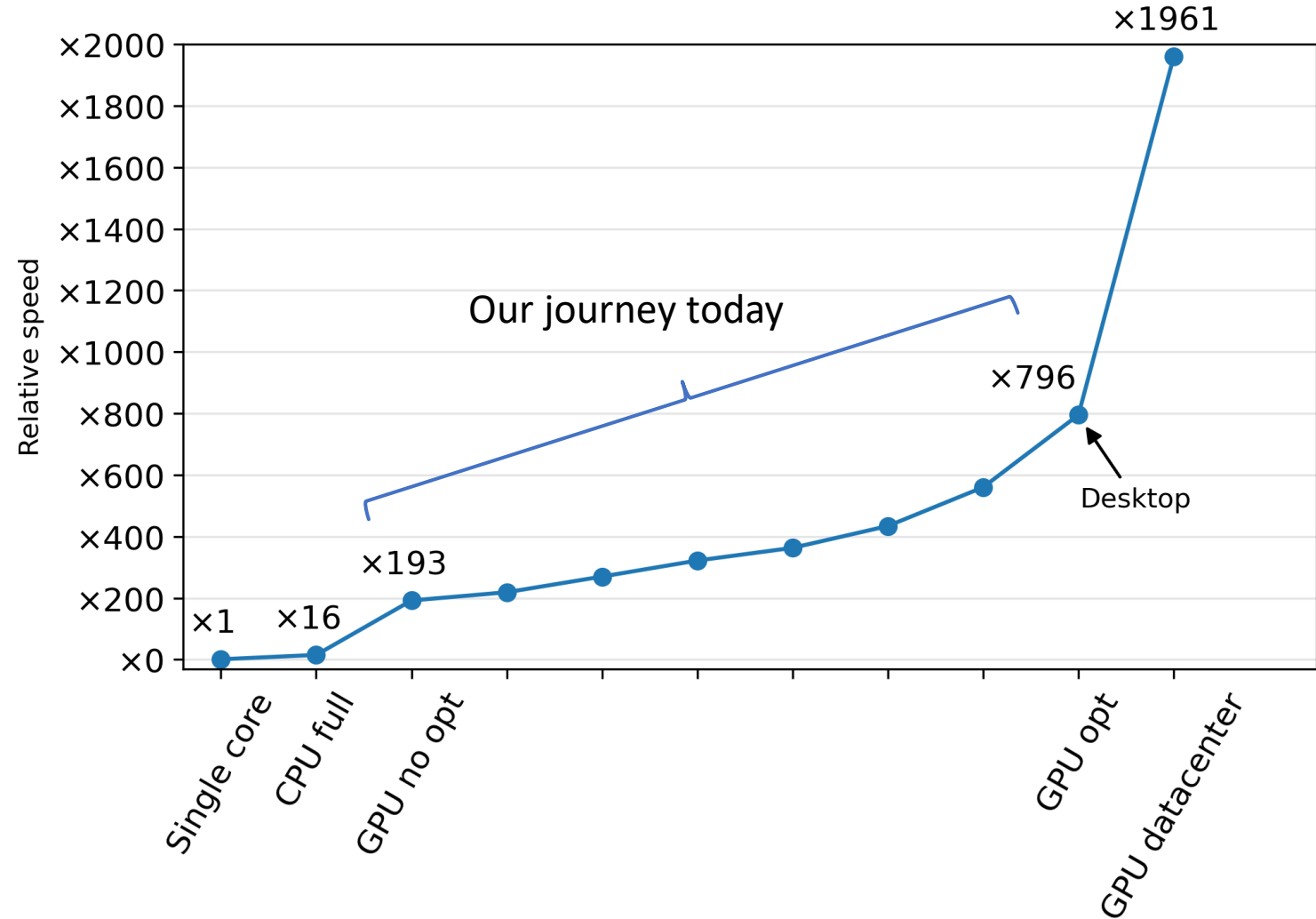
*Presented data acquired over a two years period of hands-on experimentation and accumulated frustration and suffering.

Miguel De Le Court

UCLouvain, IMMC

The importance of using GPU resources effectively

- Big gap between a naïve porting of a CPU-optimized code and a GPU-optimized code
- Writing efficient GPU code is HARD
- When it works : it's very fast



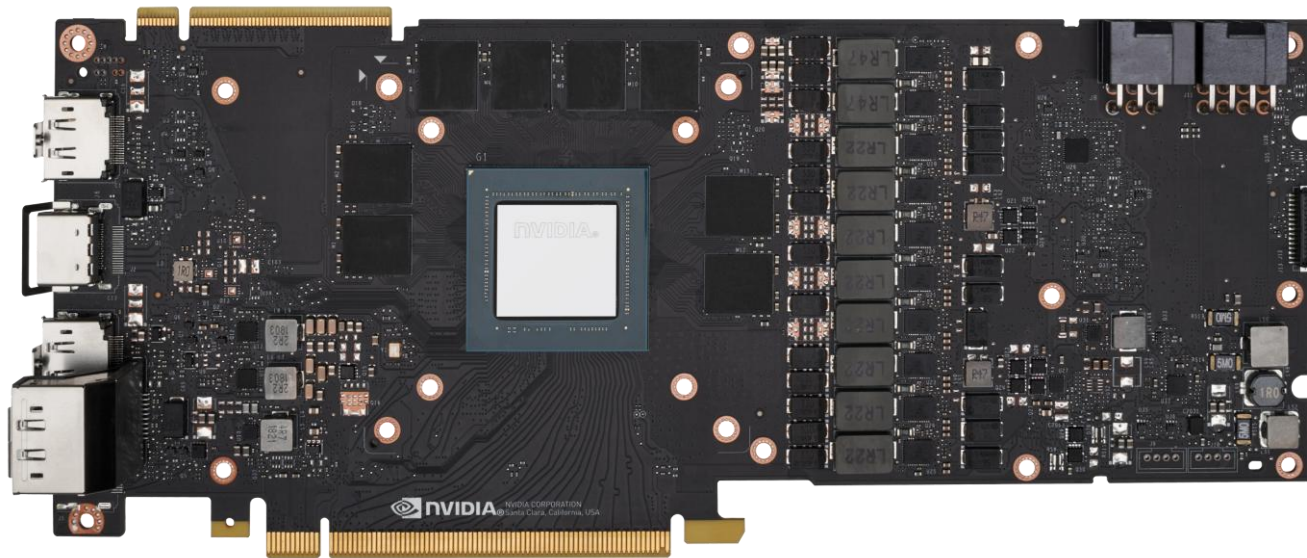
Plan for this session

0. What is a GPU
1. double vs float
2. Locality
3. Coalescence
4. double literals
5. Occupancy limiters
6. Kernel fusion
7. Shared memory
8. Array of struct of array
9. Free compiler flags

What is a GPU ?

Some differences VS a CPU include

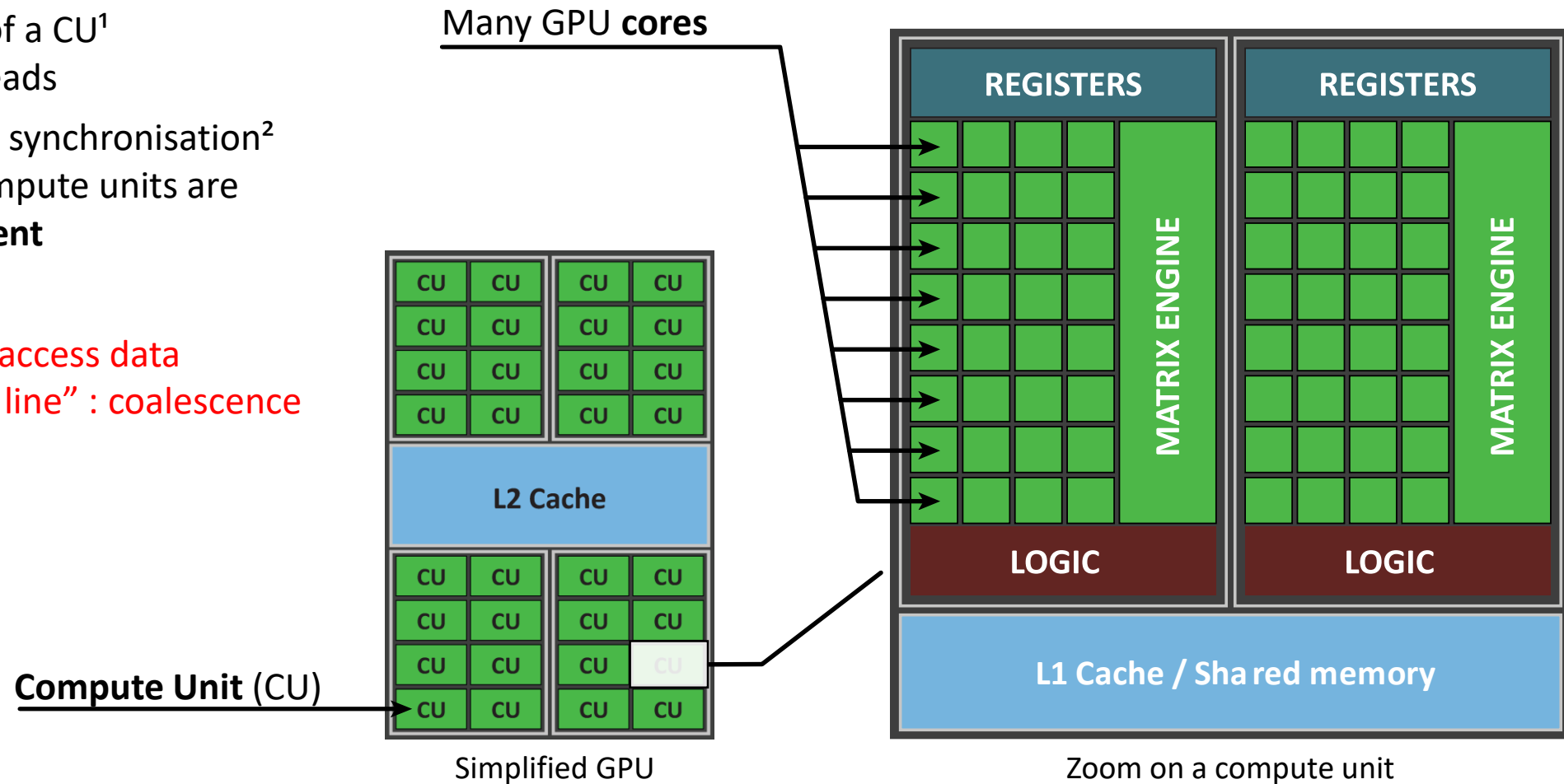
- SIMD-like execution model
- Coalescent memory access
- Very high memory latency
- Designed for higher arithmetic intensity
- Very limited cache per thread
- ...



What is a GPU : the execution model

Defining characteristics:

- Cores of a CU¹ are **not independent**
- Computations inside of a CU¹ is the same for all threads
- No data exchange and synchronisation² outside of the CU. Compute units are **completely independent**
- Cores in a CU want to access data from the same "cache line" : coalescence



Nomenclature

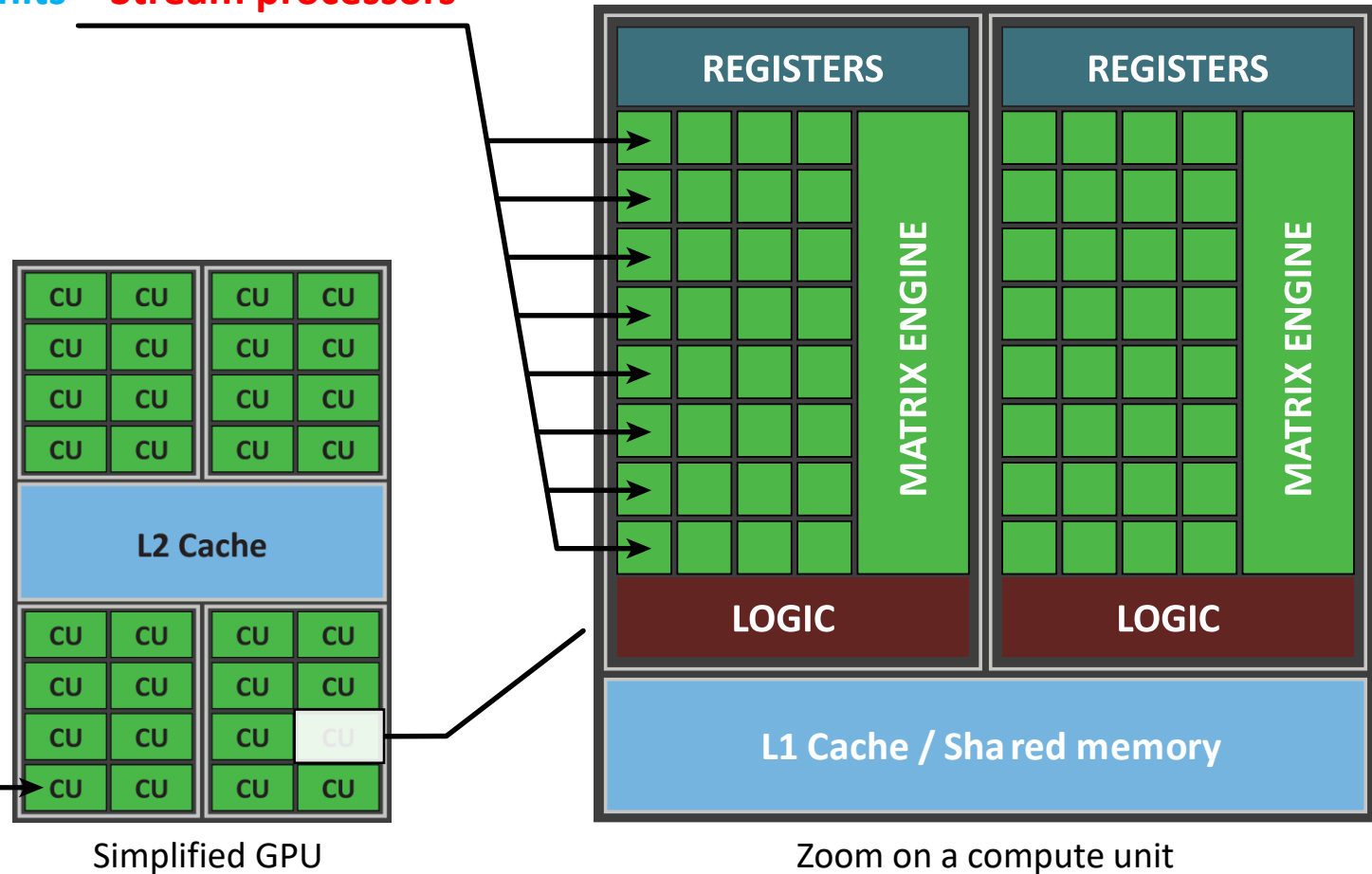
Green = NVIDIA

Blue = INTEL

Red = AMD

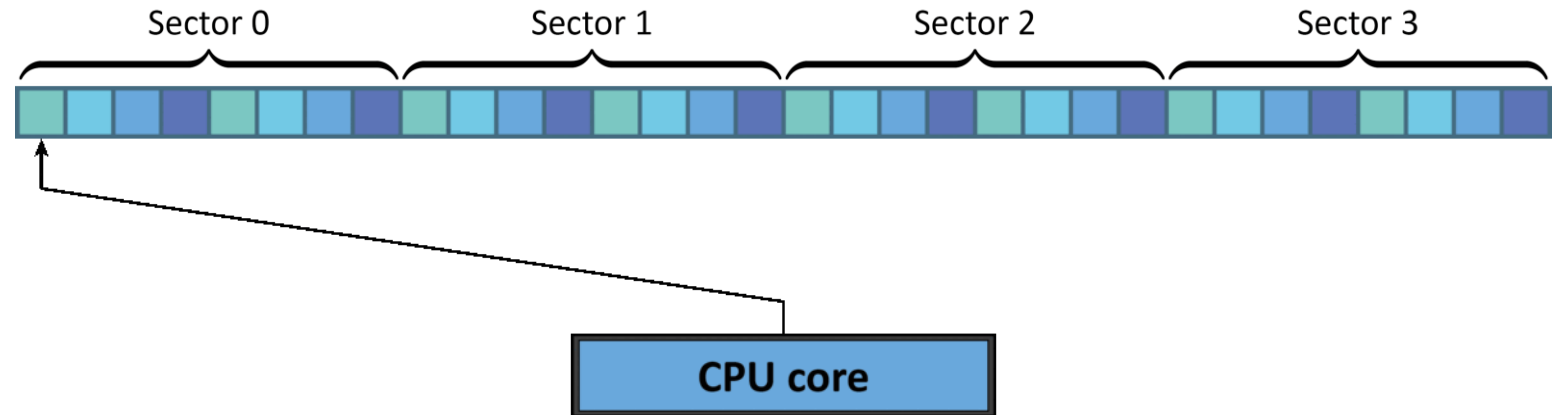
(Cuda) cores = Shading Units = Stream processors

Streaming multiprocessor = Core = Compute Unit



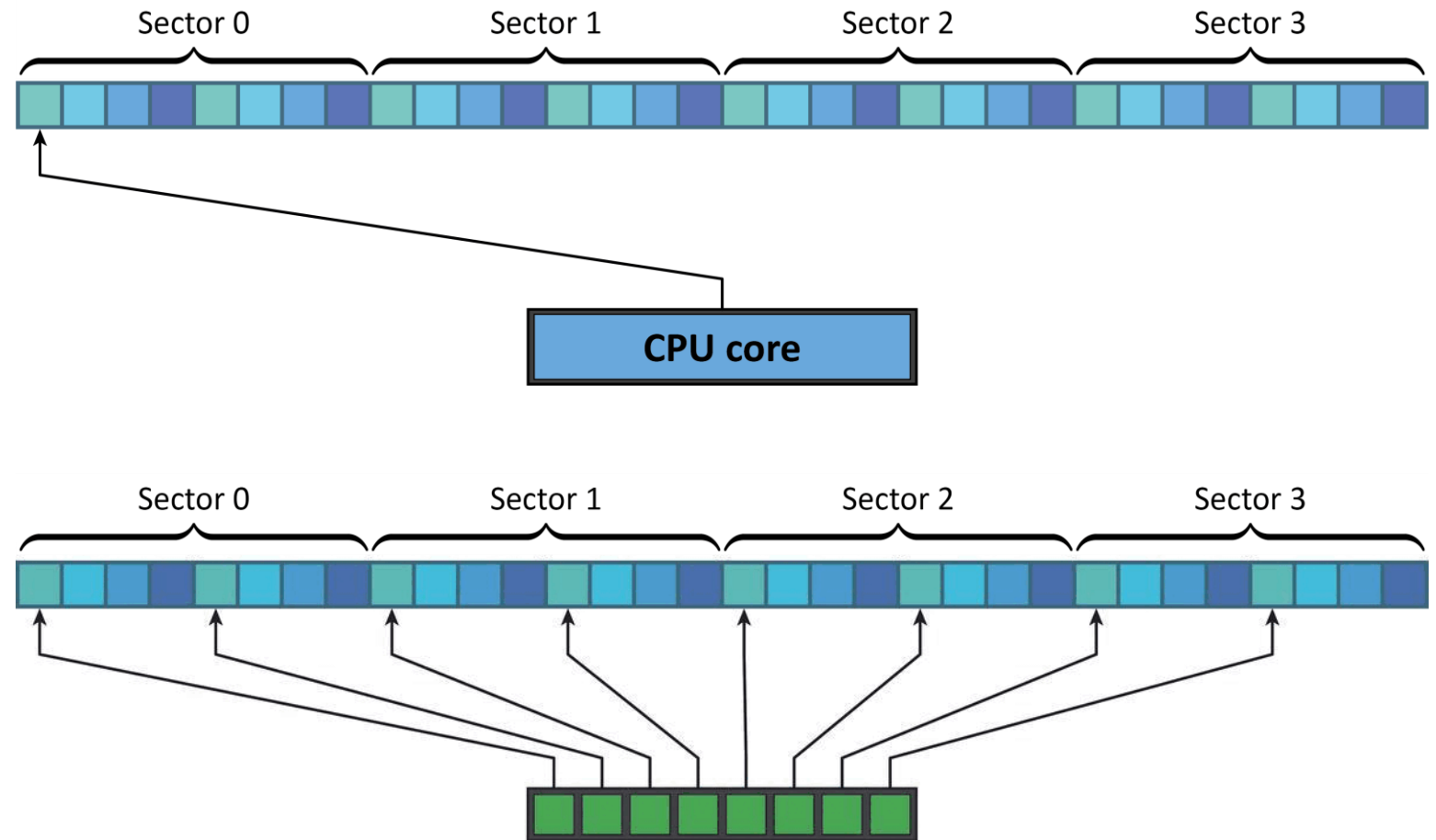
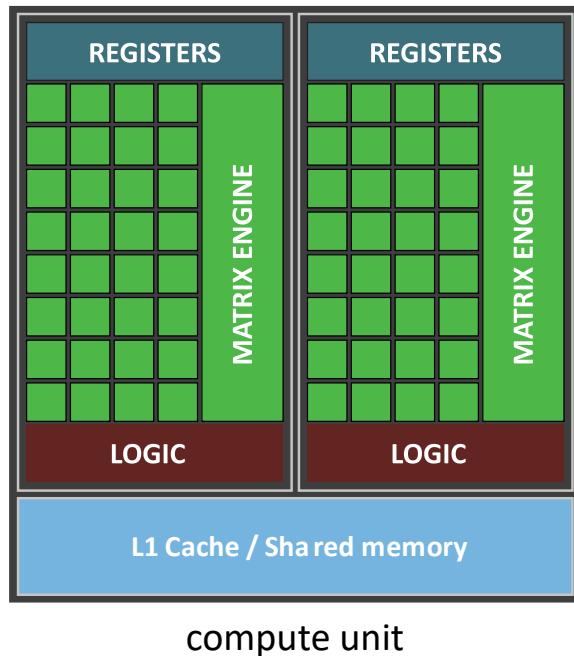
Coalescent memory access

- On the CPU : we want to maximize locality



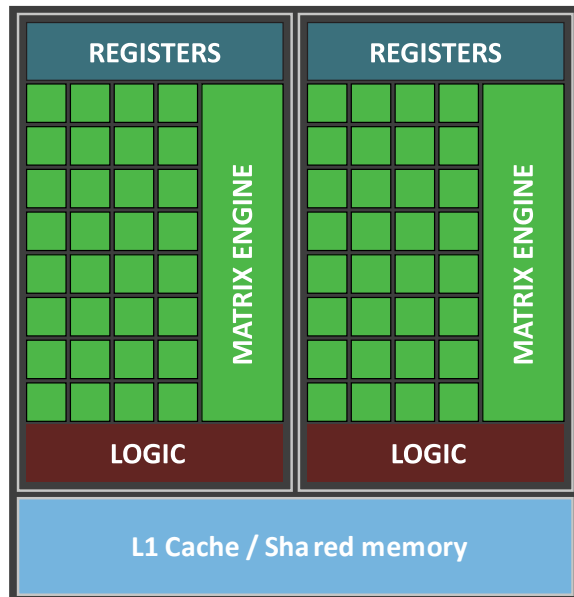
Coalescent memory access

- On the CPU : we want to maximize locality
- On the GPU : data is accessed simultaneously
- Much smaller cache per core : data may not fit! → Excessive loads

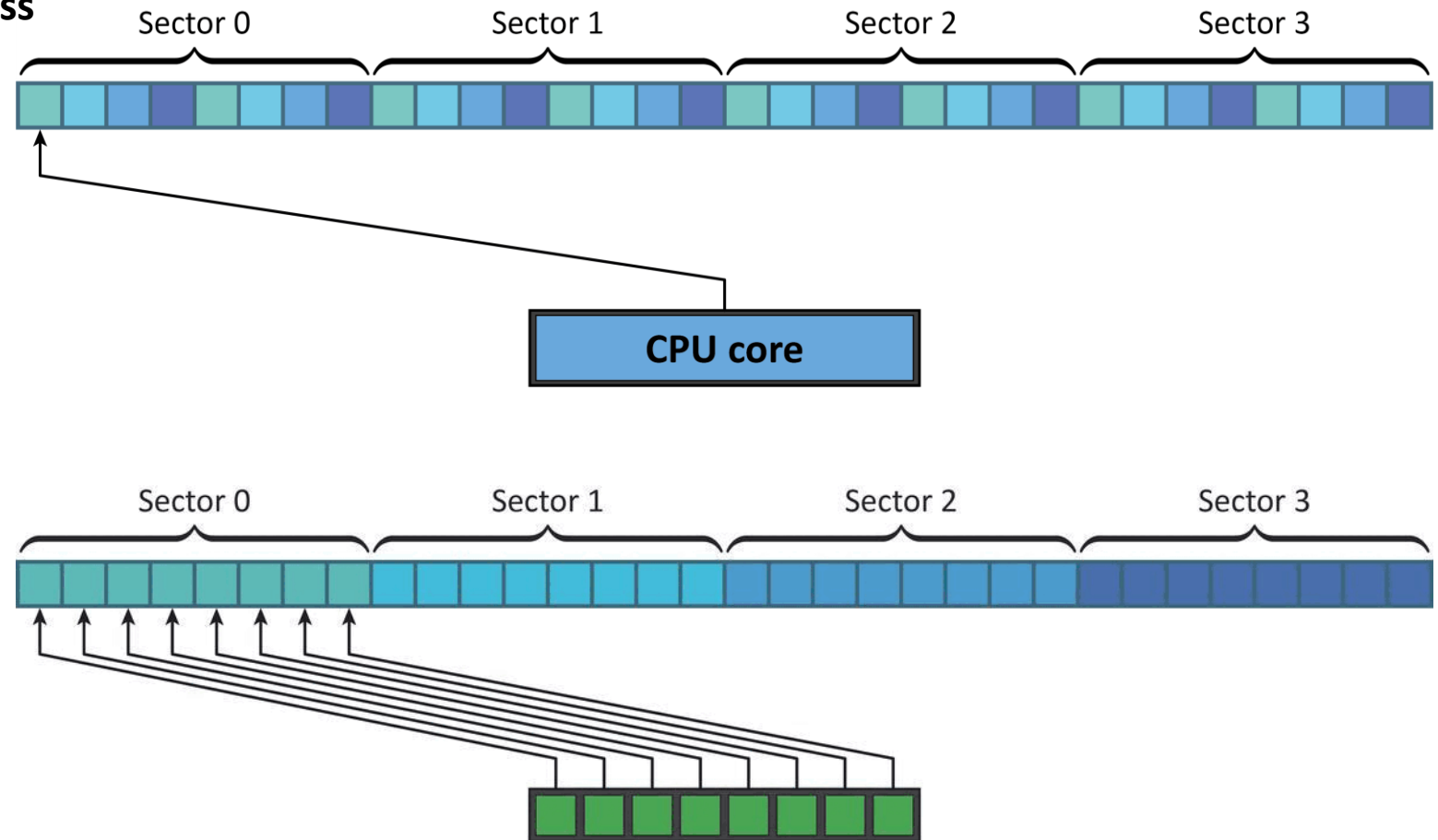


Coalescent memory access

- On the CPU : we want to maximize locality
- On the GPU : data is accessed simultaneously
- Much smaller cache per core : data may not fit! → Excessive loads
- Optimal pattern : all cores from a CU read same sector : **One sector read per access**



compute unit



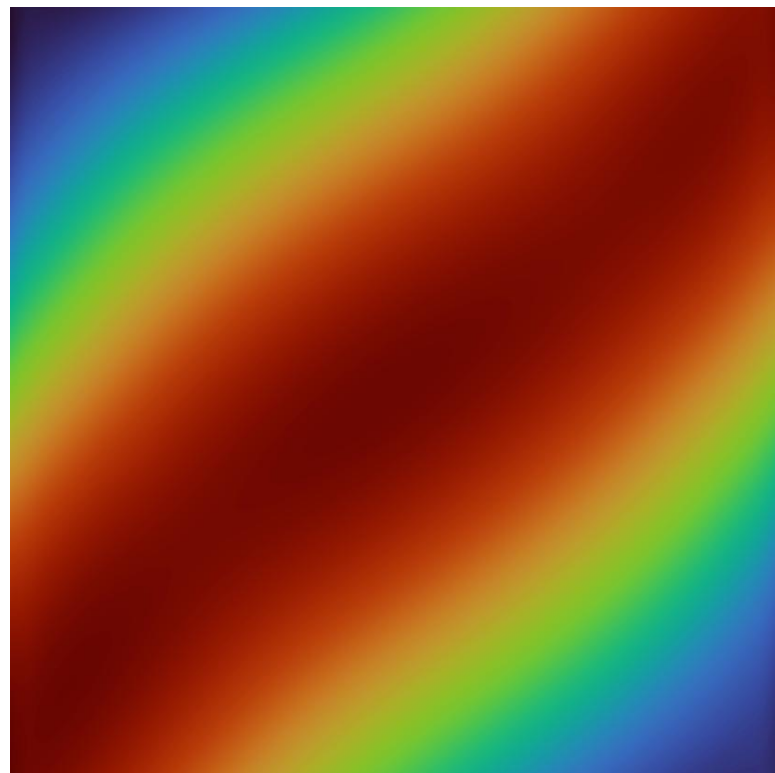
1 – Double vs float : which one should you choose, and why is it float*?

- RTX3080:

- Float : 29.77 TFLOPS
- Double : 0.47 TFLOPS

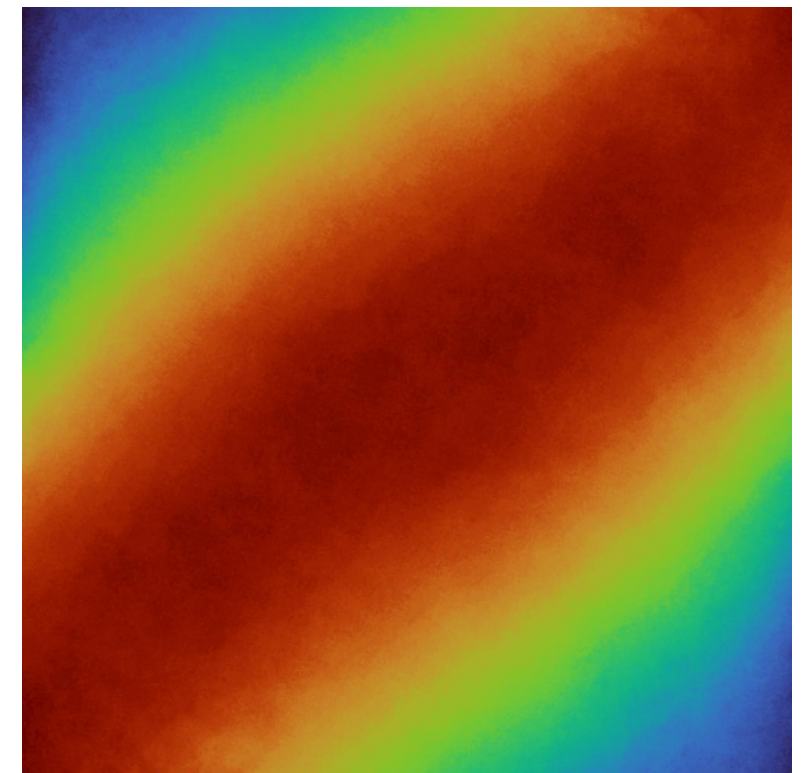
- A100:

- Float : 19.49 TFLOPS
- Double : 9.746 TFLOPS



Double

If you can, use floats



Float

Where is the bottleneck ?

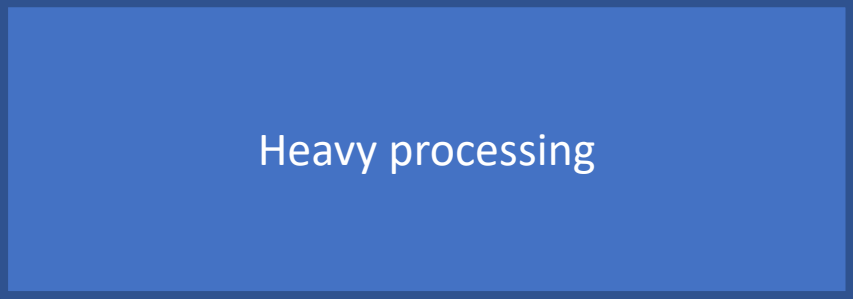
```
void dudt(const float* u, float* fu, ...){  
    for(int i = 0; i < n; i++){  
        // loading data  
        float[4][3] local_u;  
        local_u[...] = u[...];  
  
        // some heavy computation  
  
        // writing back the result  
        fu[...] = local_fu;  
    }  
}
```

The diagram illustrates the execution flow of the provided C code. A large blue rectangular box labeled "Heavy processing" is positioned between the "loading data" and "writing back the result" sections. To the right of the code, three blue brackets indicate the percentage of time spent on each section: 4% for the "loading data" section, 95% for the "Heavy processing" section, and 1% for the "writing back the result" section.

| Code Section | Percentage |
|-------------------------|------------|
| loading data | 4% |
| Heavy processing | 95% |
| writing back the result | 1% |

Where is the bottleneck ?

```
__global__ void dudt_kernel(const float* u, float* fu, ...){
    int tid = threadIdx.x + blockIdx.x * blockDim.x;
    if (tid < n){
        // loading data
        float[4][3] local_u;
        local_u[...] = u[...];

        // some heavy computation
        
        Heavy processing

        // writing back the result
        fu[...] = local_fu;
    }
}
```

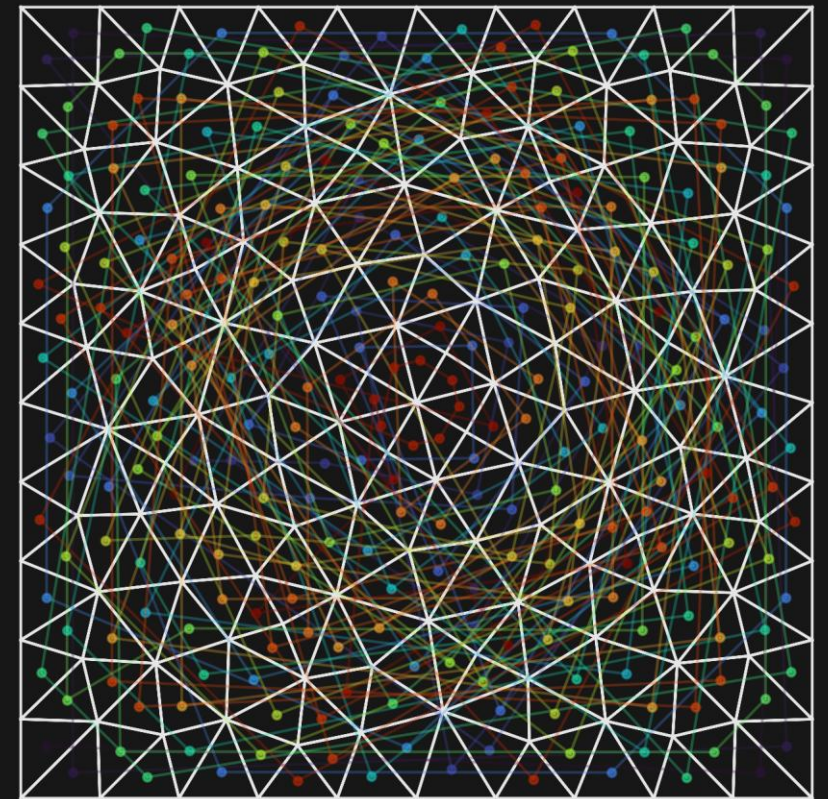
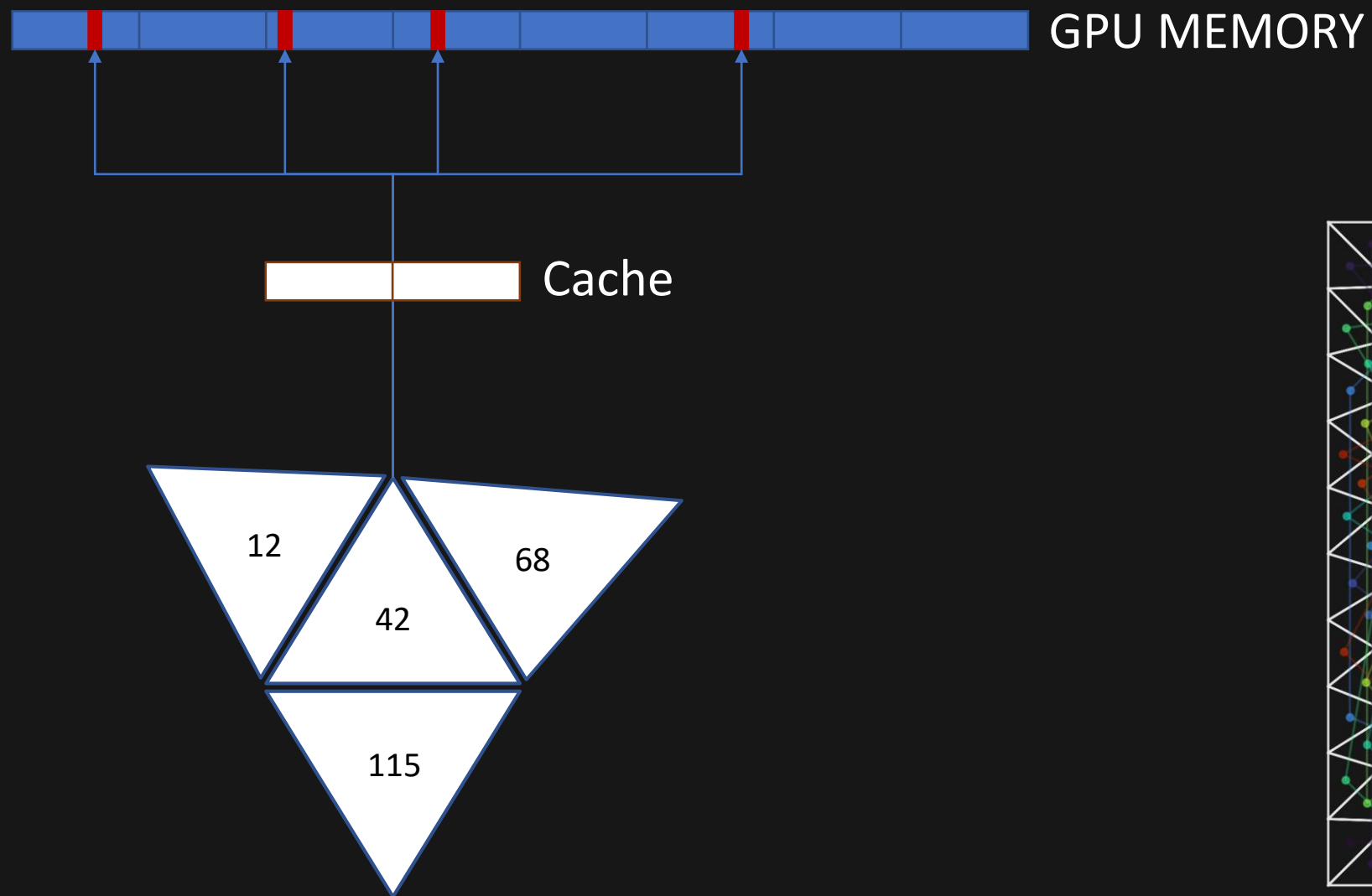
Memory is (almost always) the bottleneck

```
__global__ void dudt_kernel(const float* u, float* fu, ...){  
    int tid = threadIdx.x + blockIdx.x * blockDim.x;  
    if (tid < n){  
        // loading data  
        float[4][3] local_u;  
        local_u[...] = u[...];  
  
        // some heavy computation  
  
        // writing back the result  
        fu[...] = local_fu;  
    }  
}
```

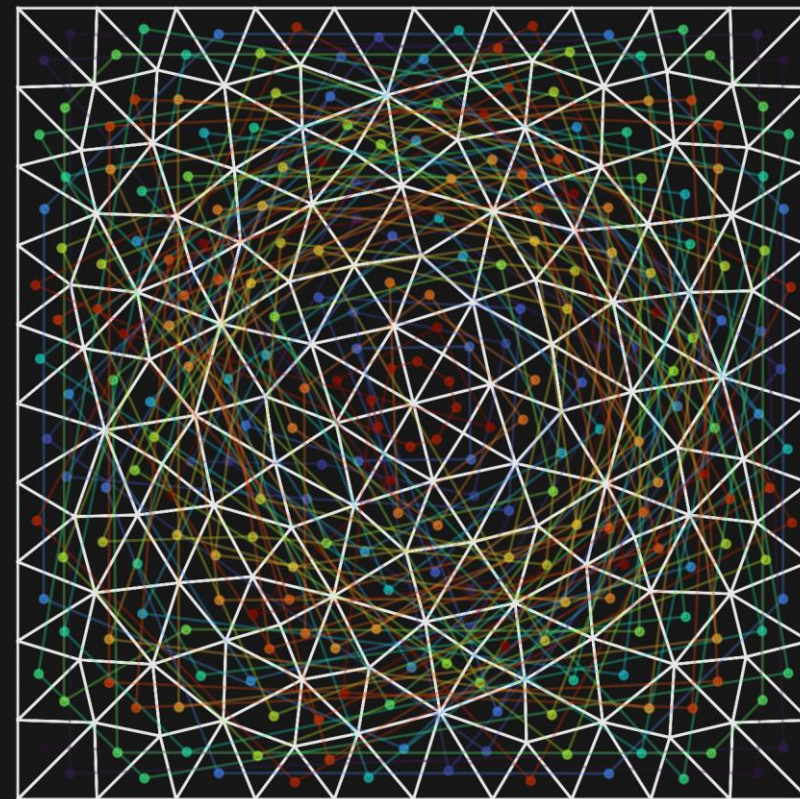
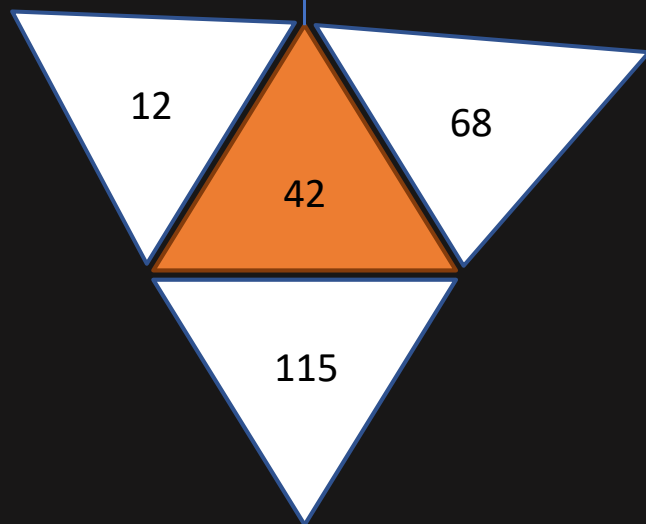
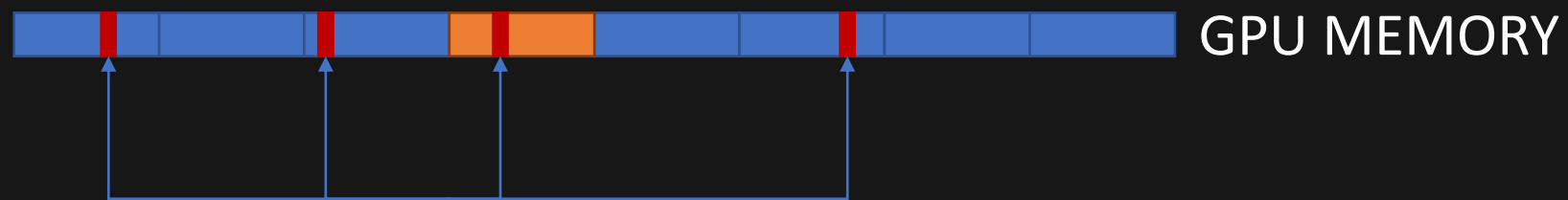
The diagram illustrates the execution flow of the provided code. It is divided into three main sections, each with a percentage value indicated by a bracket on the right side:

- 55%**: This section covers the code from `local_u[...] = u[...];` to the start of the heavy computation. It is annotated with the comment `// loading data`.
- 30%**: This section is represented by a large blue rectangular box containing the text "Heavy processing". It is annotated with the comment `// some heavy computation`.
- 15%**: This section covers the code from `fu[...] = local_fu;` to the end of the function. It is annotated with the comment `// writing back the result`.

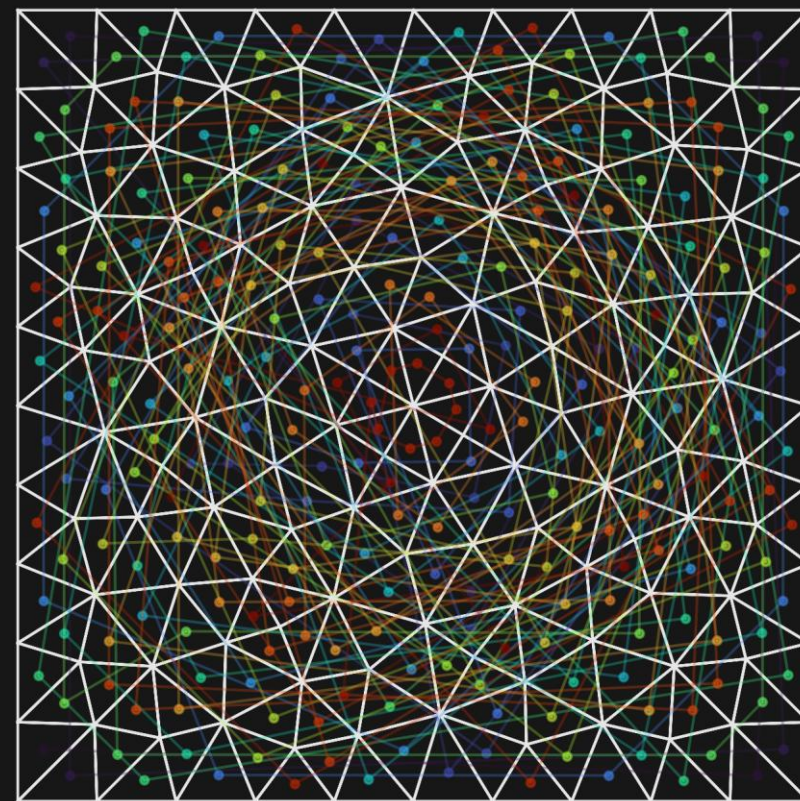
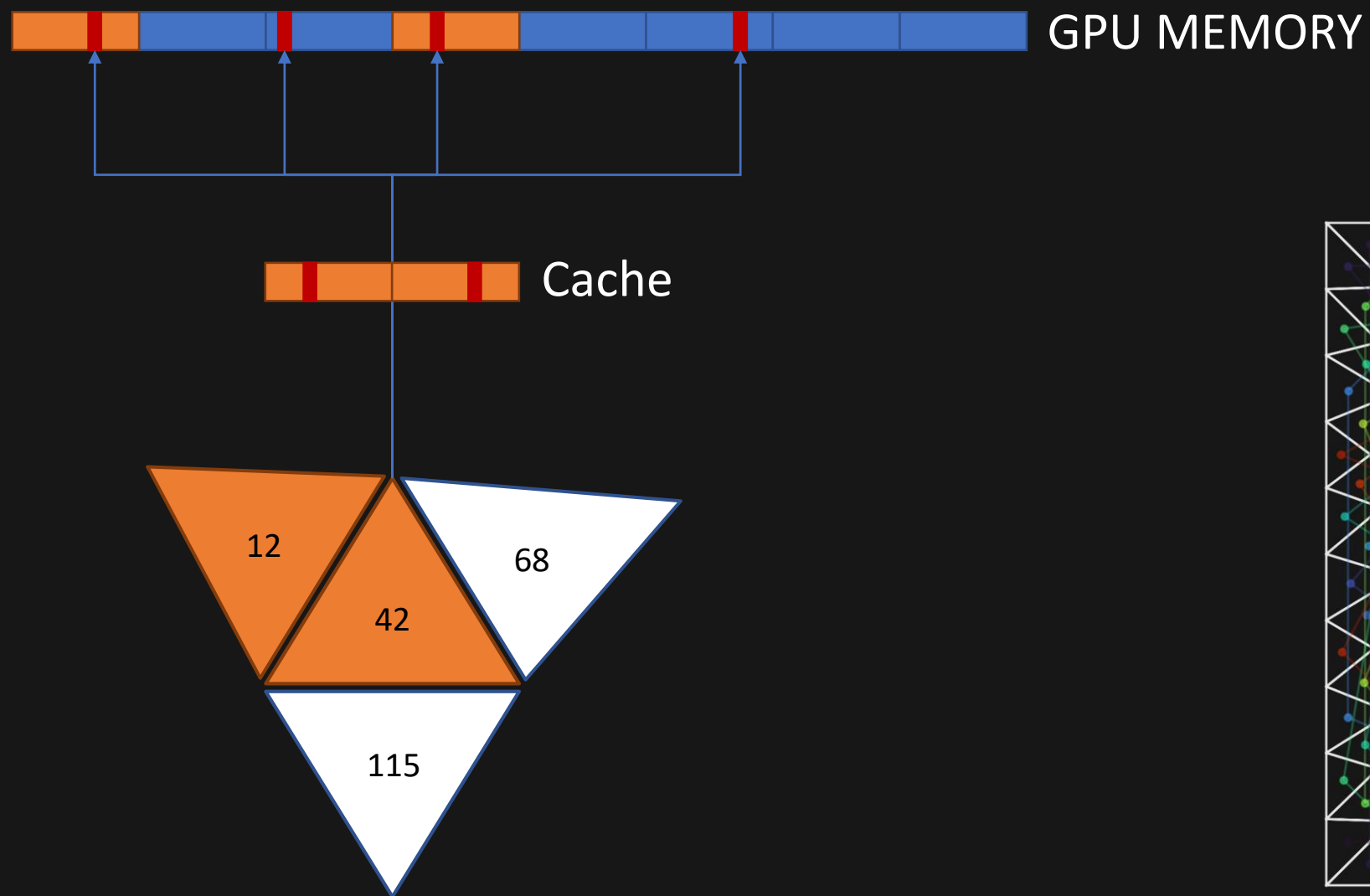
2 – Importance of memory locality



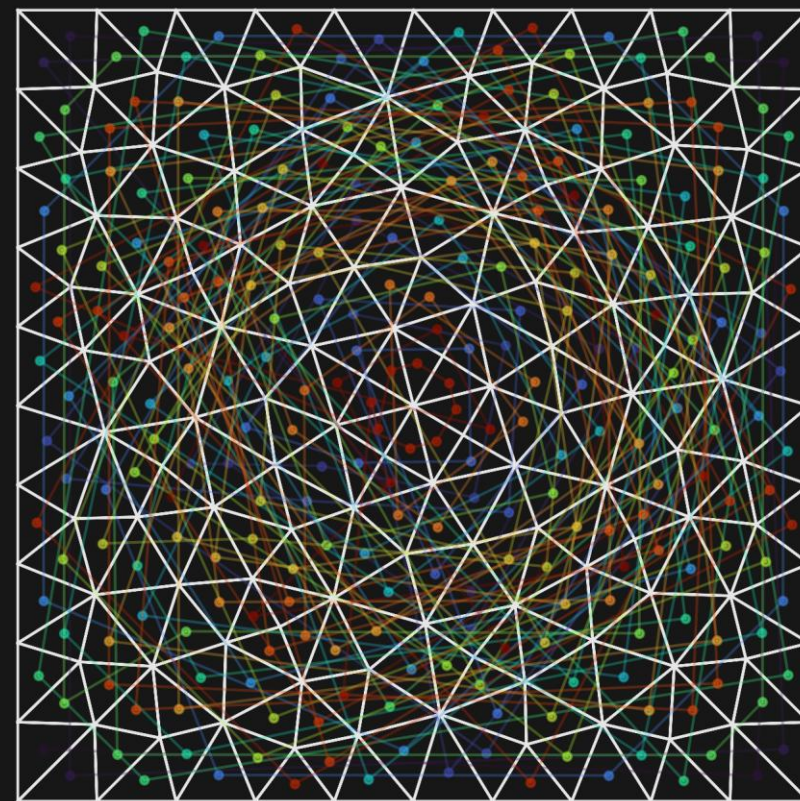
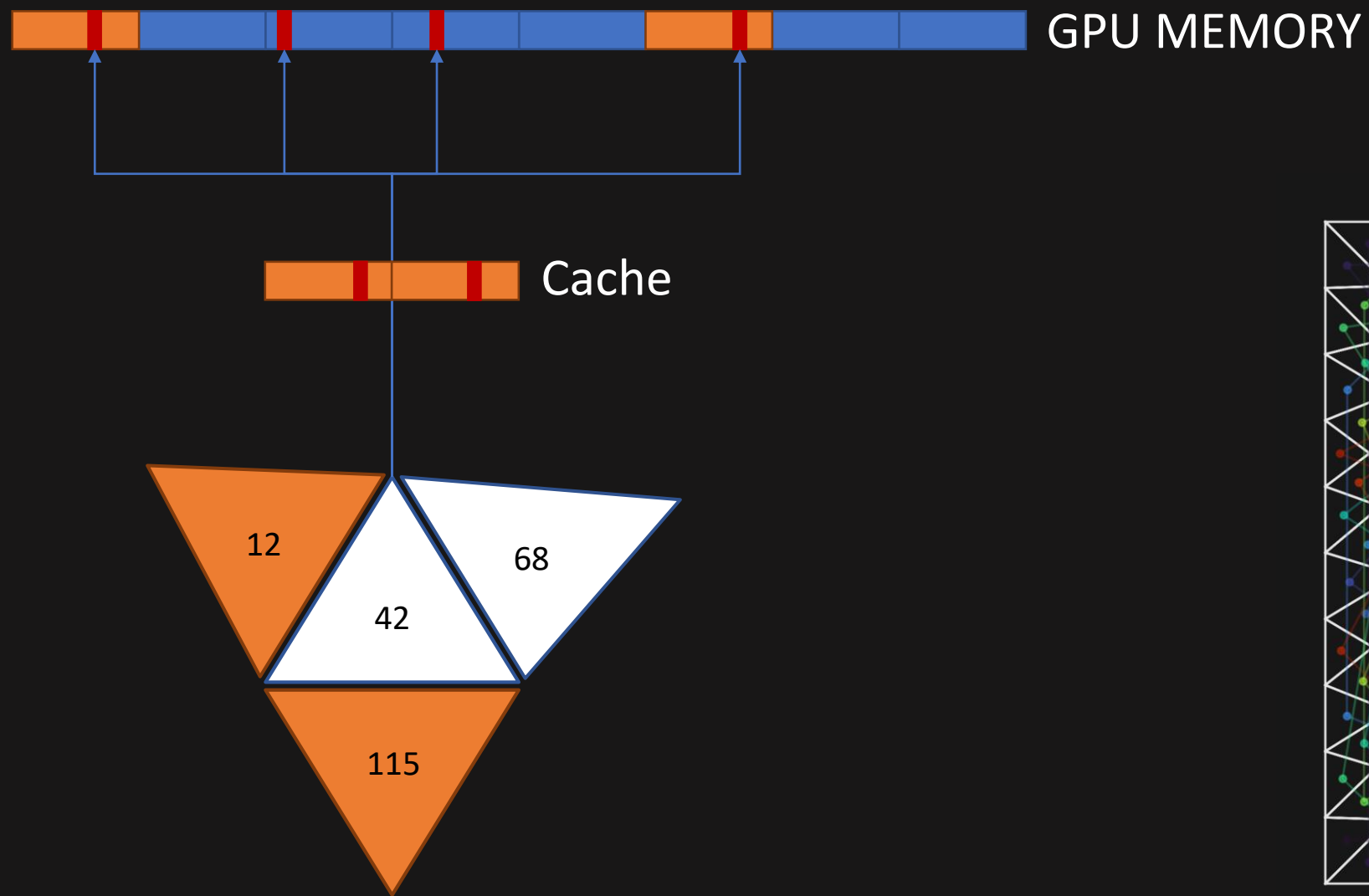
2 – Importance of memory locality



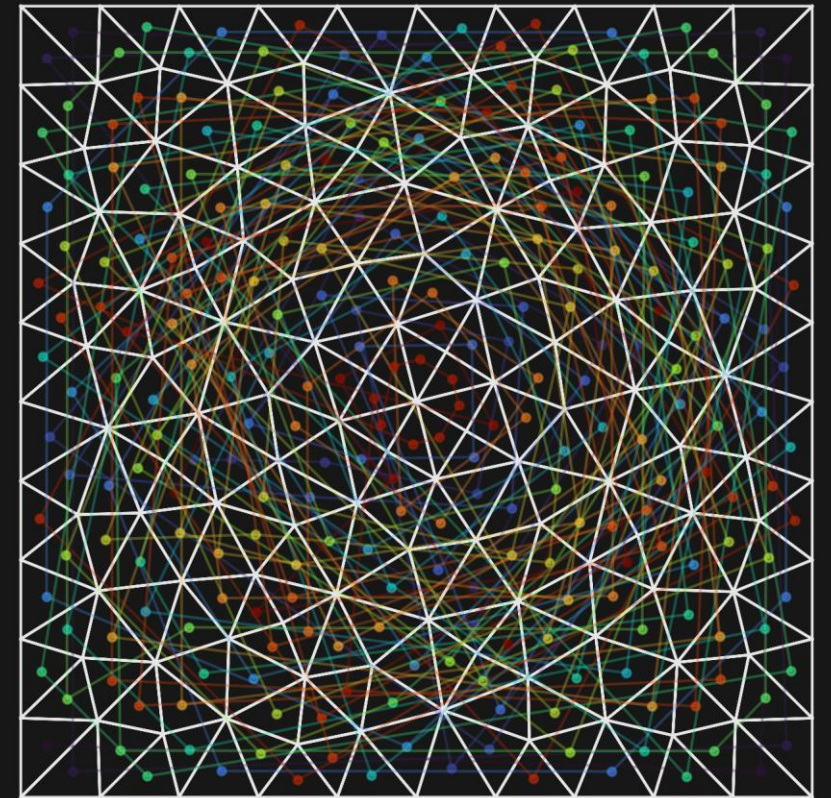
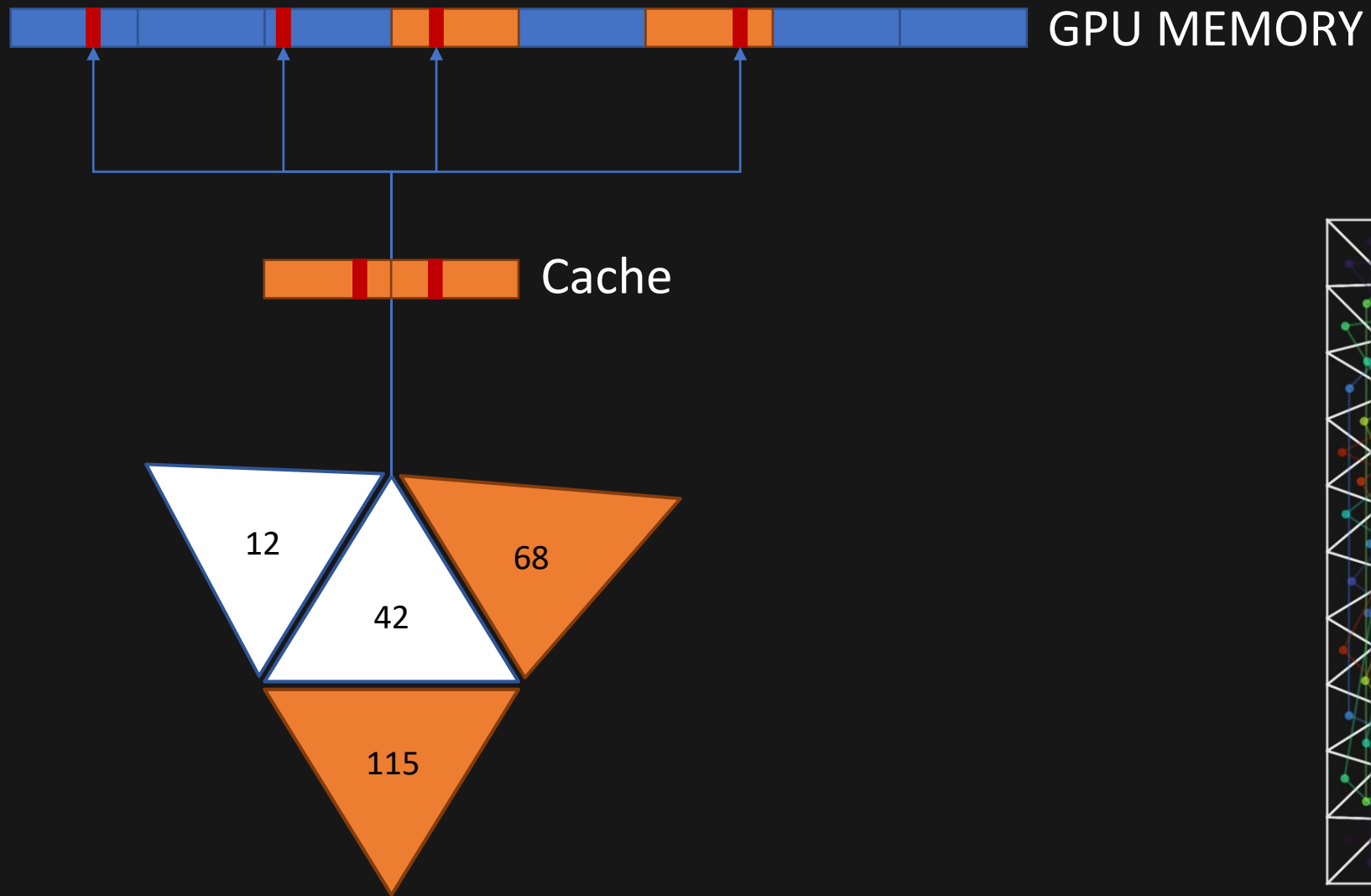
2 – Importance of memory locality



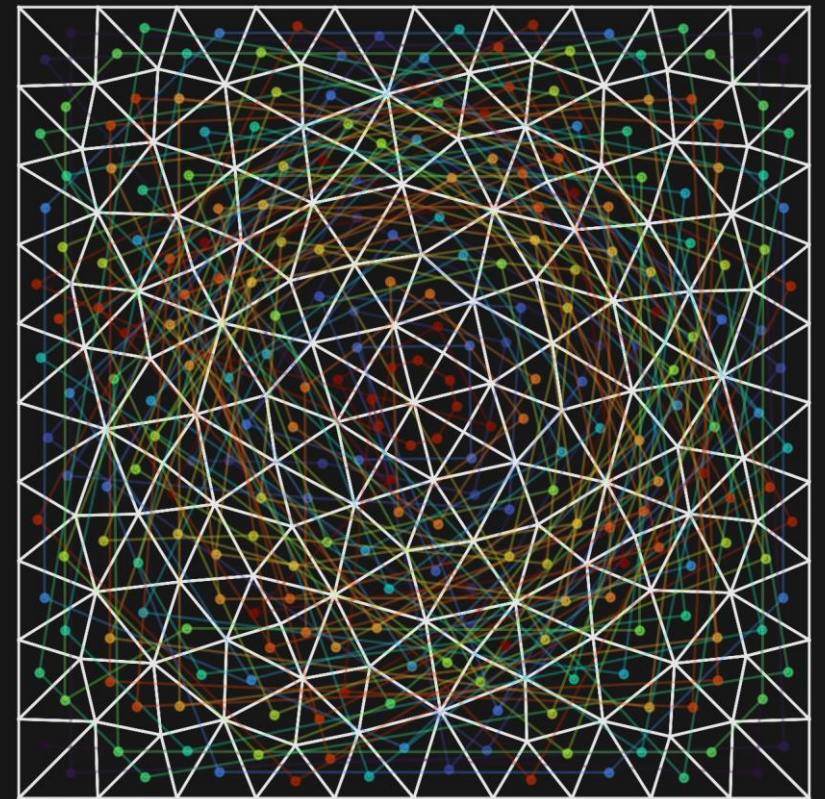
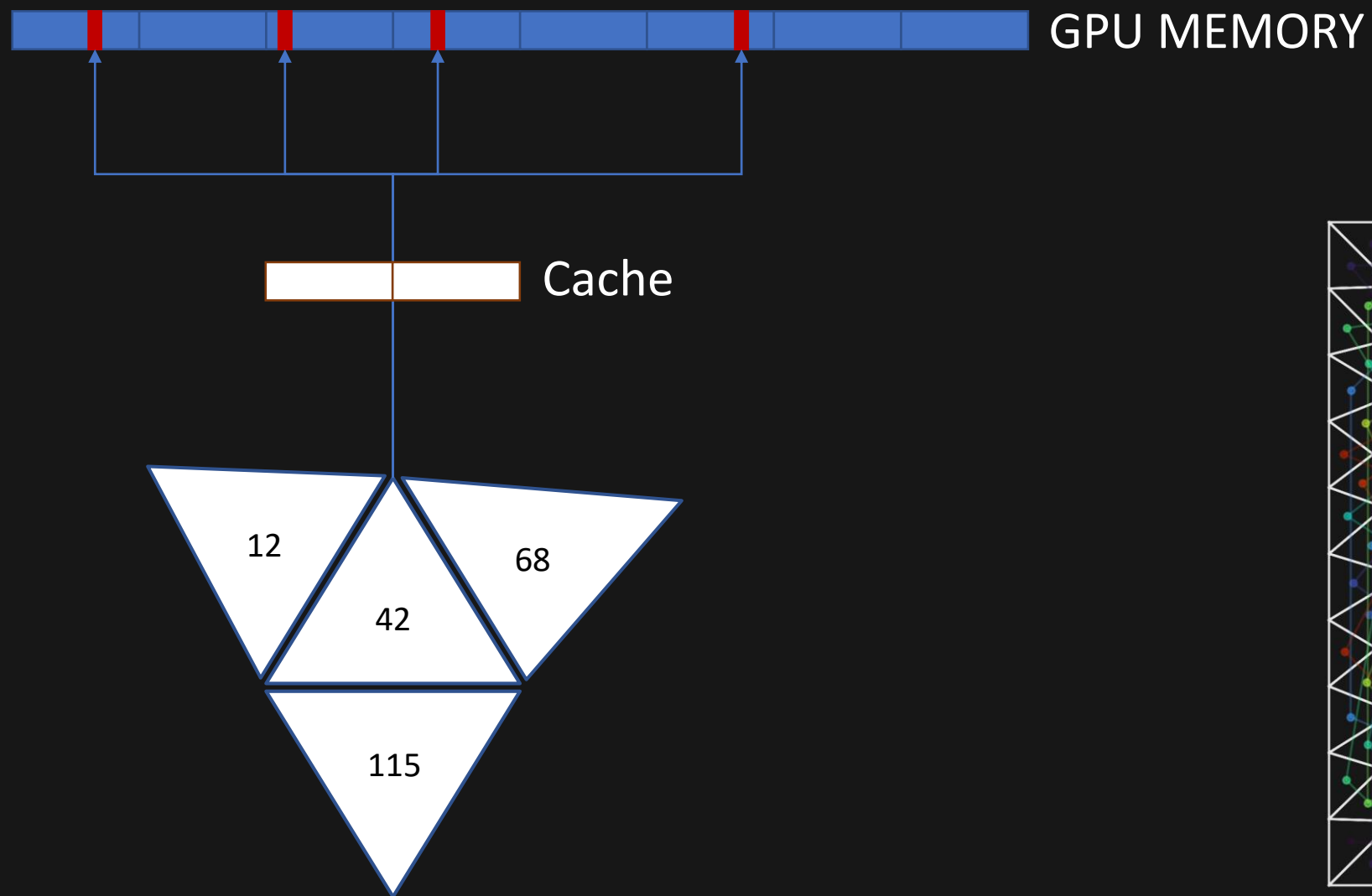
2 – Importance of memory locality



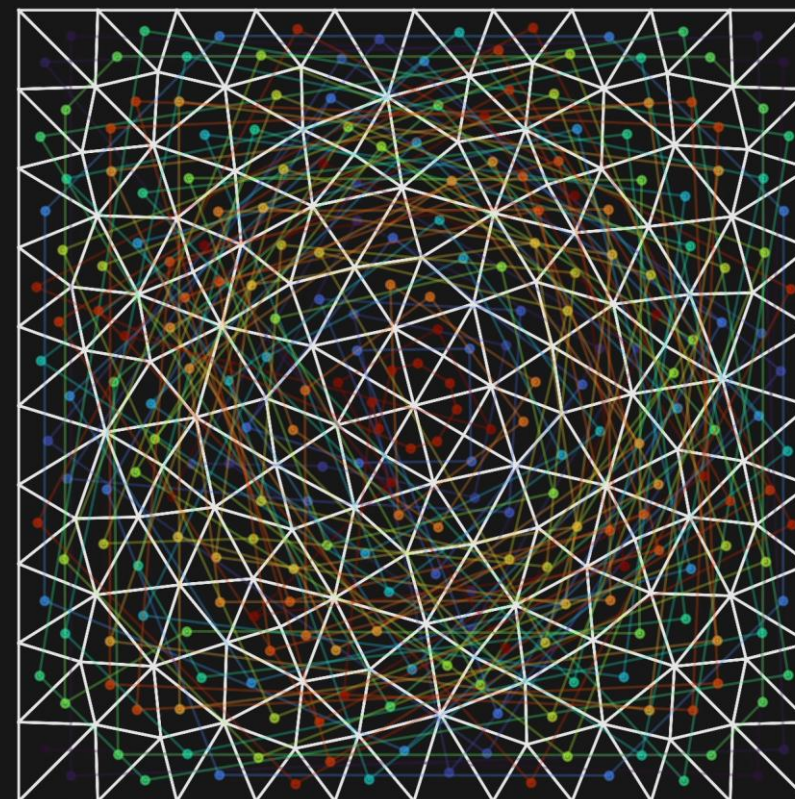
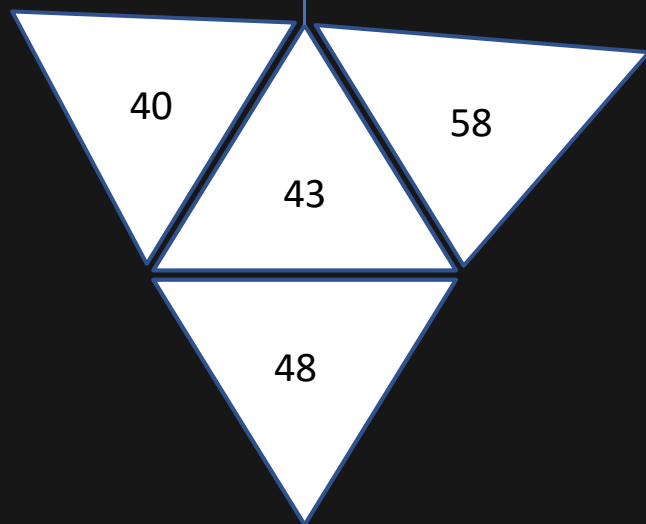
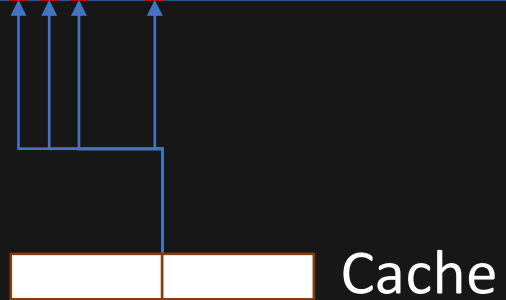
2 – Importance of memory locality



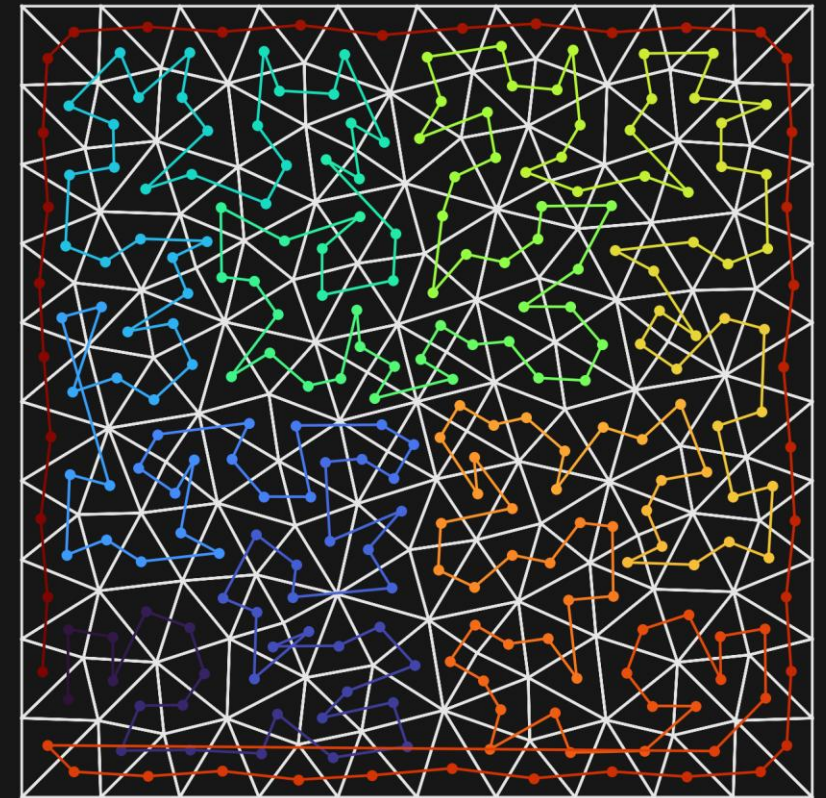
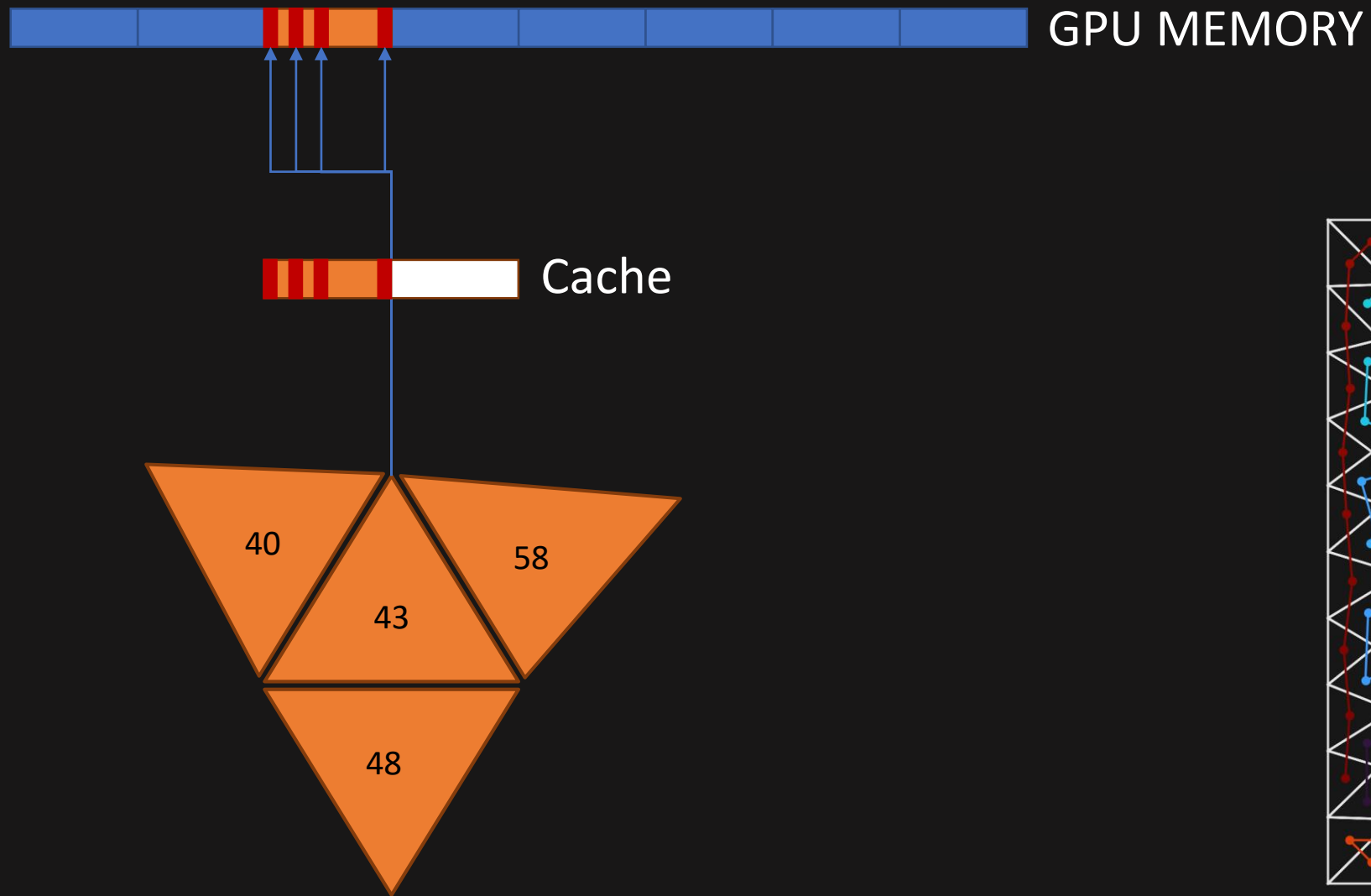
2 – Importance of memory locality



2 – Importance of memory locality



2 – Importance of memory locality



2 – Clues to detect that problem:

- Very high bandwidth utilization
- Too many reads compared to what is expected

L2 Load Access Pattern The memory access pattern for loads from L1TEX to L2 is not optimal. The granularity of an L1TEX request to L2 is a 128 byte cache line. That is 4 consecutive 32-byte sectors per L2 request. However, this kernel only accesses an average of 1.2 sectors out of the possible 4 sectors per cache line. Check the [Source Counters](#) section for uncoalesced loads and try to minimize how many cache lines need to be accessed per memory request.

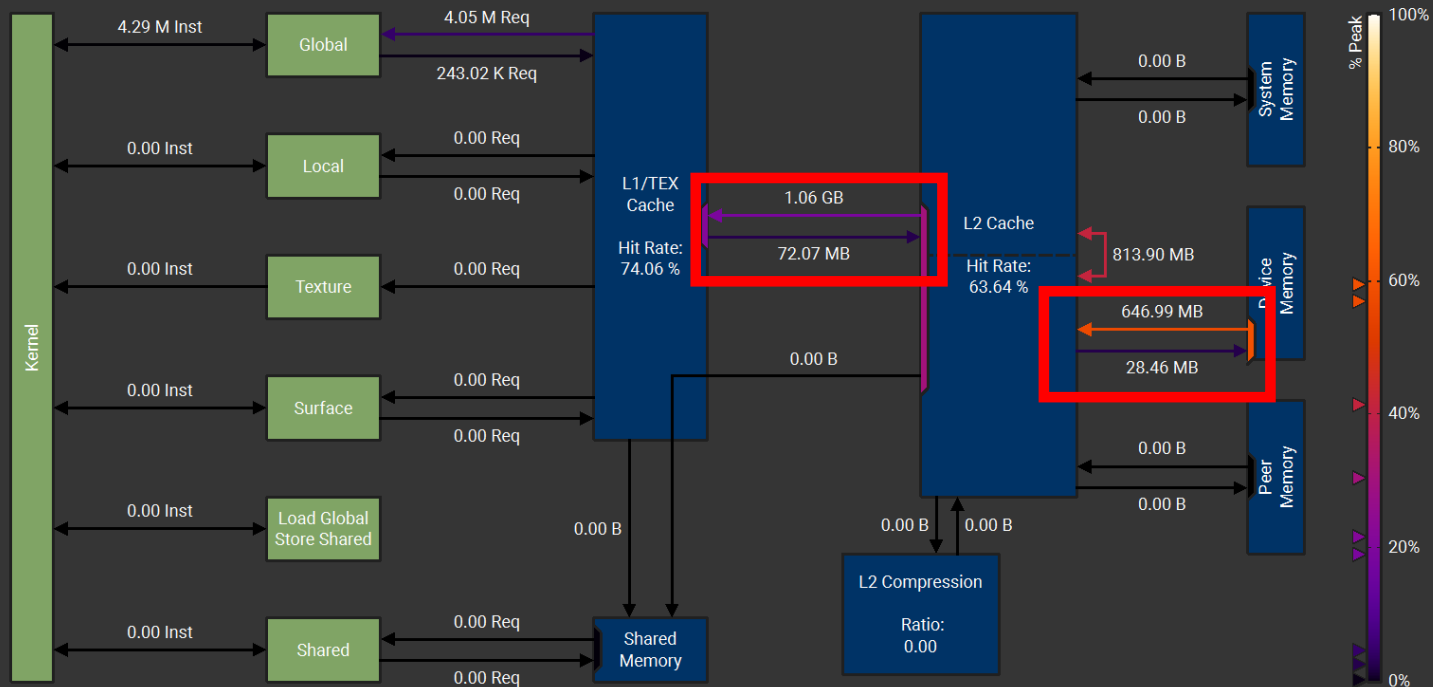
DRAM Excessive Read Sectors The memory access pattern for loads from device memory causes 21,733,744 sectors to be read from DRAM, which is 1.8x of the 12,160,697 sectors causing a miss in the L2 cache. The DRAM fetch granularity for read misses in L2 is 64 bytes, i.e. the lower or upper half of an L2 cache line. Try changing your access pattern to make use of both sectors returned by a DRAM read request for optimal usage of the DRAM throughput. For strided memory reads, avoid strides of 64 bytes or larger to avoid moving unused sectors from DRAM to L2.

Shared Memory Conflicts Detection of shared memory bank conflicts.

Apply

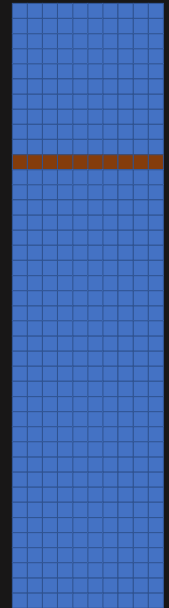
Memory Chart

Show As: Transfer Size

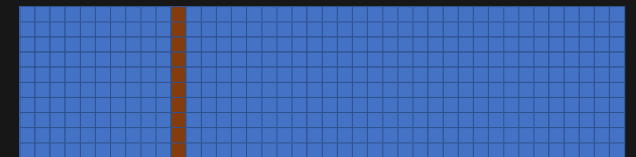


3 – So... let's optimize for locality, right ?

```
__global__ void mysum(const float* a, const float* b, float* c, int nf, int n){
    int tid = threadIdx.x + blockIdx.x * blockDim.x;
    if (tid < n){
        for(int field = 0; field < nf; field++){
            // maximal locality
            c[tid * nf + field] = a[tid * nf + field] + b[tid * nf + field];
        }
    }
}
```



```
__global__ void mysum(const float* a, const float* b, float* c, int nf, int n){
    int tid = threadIdx.x + blockIdx.x * blockDim.x;
    if (tid < n){
        for(int field = 0; field < nf; field++){
            // very large stride
            c[field * n + tid] = a[field * n + tid] + b[field * n + tid];
        }
    }
}
```



3 – No, we need to consider coalescence

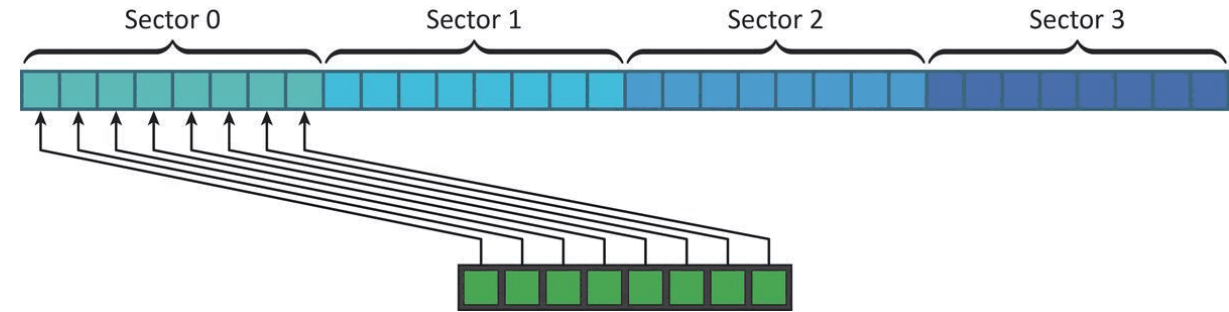
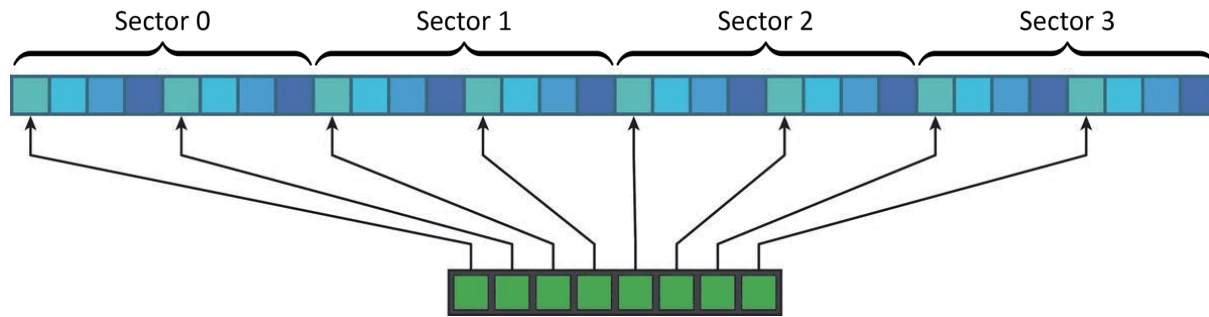
```
__global__ void mysum(const float* a, const float* b, float* c, int nf, int n){  
    int tid = threadIdx.x + blockIdx.x * blockDim.x;  
    if (tid < n){  
        for(int field = 0; field < nf; field++){  
            // maximal locality  
            c[tid * nf + field] = a[tid * nf + field] + b[tid * nf + field];  
        }  
    }  
}
```

634 μ s

```
__global__ void mysum(const float* a, const float* b, float* c, int nf, int n){  
    int tid = threadIdx.x + blockIdx.x * blockDim.x;  
    if (tid < n){  
        for(int field = 0; field < nf; field++){  
            // very large stride  
            c[field * n + tid] = a[field * n + tid] + b[field * n + tid];  
        }  
    }  
}
```

403 μ s

3 – No, we need to consider coalescence

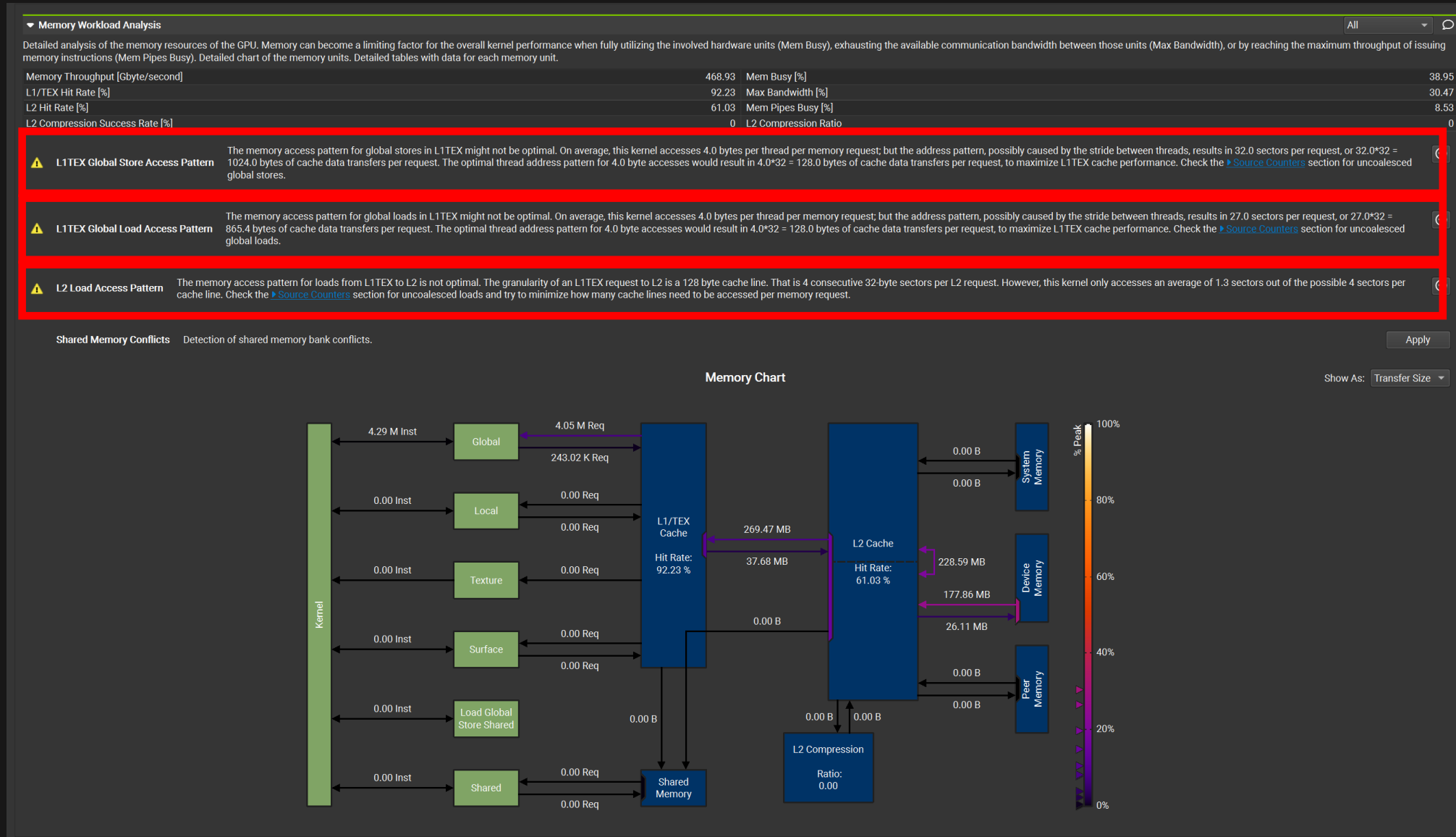


```
__global__ void mysum(const float* a, const float* b, float* c, int nf, int n){  
    int tid = threadIdx.x + blockIdx.x * blockDim.x;  
    if (tid < n){  
        for(int field = 0; field < nf; field++){  
            // maximal locality  
            c[tid * nf + field] = a[tid * nf + field] + b[tid * nf + field];  
        }  
    }  
}
```

```
__global__ void mysum(const float* a, const float* b, float* c, int nf, int n){  
    int tid = threadIdx.x + blockIdx.x * blockDim.x;  
    if (tid < n){  
        for(int field = 0; field < nf; field++){  
            // very large stride  
            c[field * n + tid] = a[field * n + tid] + b[field * n + tid];  
        }  
    }  
}
```

3 – How to evaluate if coalescence is good ?

- ncu will yell at you if it's not



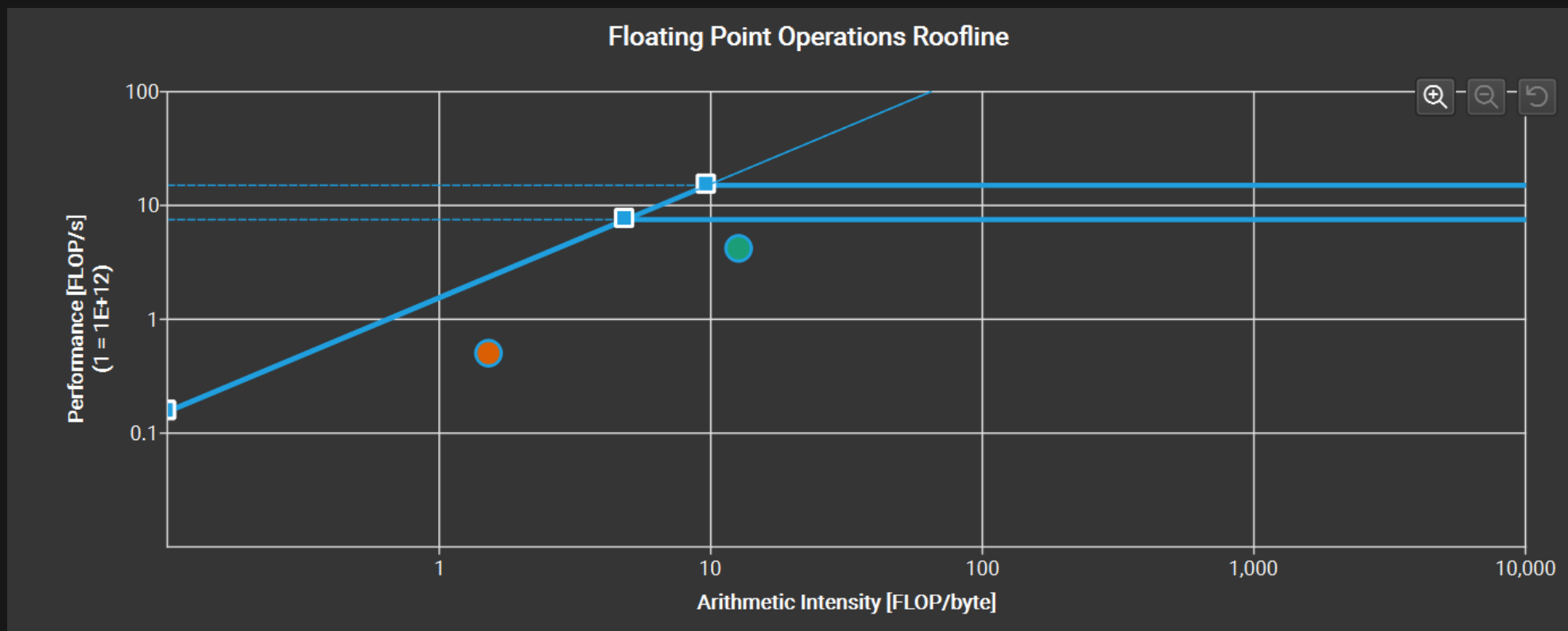
4 – Beware of doubles

```
__device__ float inv2x2(const float mat[2][2], float
inv[2][2])
{
    float det = det2x2(mat);
    float ud = 1.0 / det;
    inv[0][0] = mat[1][1] * ud;
    inv[1][0] = -mat[1][0] * ud;
    inv[0][1] = -mat[0][1] * ud;
    inv[1][1] = mat[0][0] * ud;
    return det;
}
```

4 – Beware of doubles

```
__device__ float inv2x2(const float mat[2][2], float
inv[2][2])
{
    float det = det2x2(mat);
    float ud = 1.0f / det;
    inv[0][0] = mat[1][1] * ud;
    inv[1][0] = -mat[1][0] * ud;
    inv[0][1] = -mat[0][1] * ud;
    inv[1][1] = mat[0][0] * ud;
    return det;
}
```

4 – Beware of doubles : how to check?



4 – Beware of doubles : how to check?

- Look for DMUL, DFMA
- Code must be compiled with `--generate-line-info`

The screenshot displays the NVIDIA Nsight Systems interface. At the top, the 'Page' dropdown is set to 'Source' and the 'Result' is '0 - 520 - dudt'. The 'View' is set to 'Source and SASS'. The 'Source' is 'main.h' and the 'Navigation' is 'Instructions Executed'. The 'DFMA' instruction is highlighted in the assembly view, and the corresponding source code line is highlighted in the source view.

Source code (left panel):

```
... scalar un = (unl+unr)/2.0;
... flux_t += un*(un > 0 ? Hutl : Hutr);
... flux_n += 0.5*(Hunl*unl + Hunr*unr);
```

Assembly instructions (right panel):

```
406 00007fe9 62fbb650 DFMA R4, R4, 0.5, R12
```

Summary table (right panel):

| # | Address | Source | Live Registers | arp | Stall | Sampl (All Cycl) |
|-----|-------------------|-----------------------|----------------|-------|-------|------------------|
| 406 | 00007fe9 62fbb650 | DFMA R4, R4, 0.5, R12 | 48 | 0.03% | | |

Select a source line in an active inline function to show additional information.

5 – Why did a simple `printf` make my code 1.2x slower ?

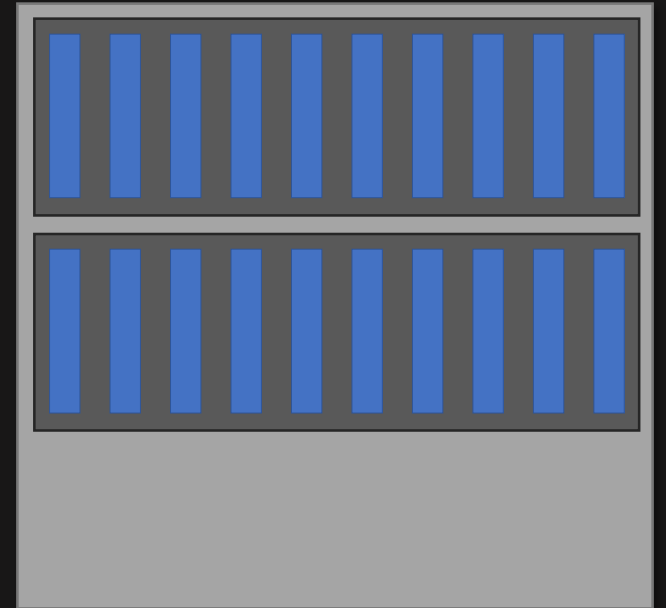
```
__device__ float inv2x2(const float mat[2][2], float
inv[2][2])
{
    float det = det2x2(mat);
    if(det){
        float ud = 1.0 / det;
        inv[0][0] = mat[1][1] * ud;
        inv[1][0] = -mat[1][0] * ud;
        inv[0][1] = -mat[0][1] * ud;
        inv[1][1] = mat[0][0] * ud;
    }
    else{
        printf("Singular matrix 2x2");
        for(int i = 0; i < 2; i++)
            for(int j = 0; j < 2; j++)
                inv[i][j] = 0.0;
    }
    return det;
}
```

5 – Why did a simple `printf` make my code 1.2x slower ?

```
__device__ float inv2x2(const float mat[2][2], float
inv[2][2])
{
    float det = det2x2(mat);
    float ud = 1.0 / det;
    inv[0][0] = mat[1][1] * ud;
    inv[1][0] = -mat[1][0] * ud;
    inv[0][1] = -mat[0][1] * ud;
    inv[1][1] = mat[0][0] * ud;
    return det;
}
```

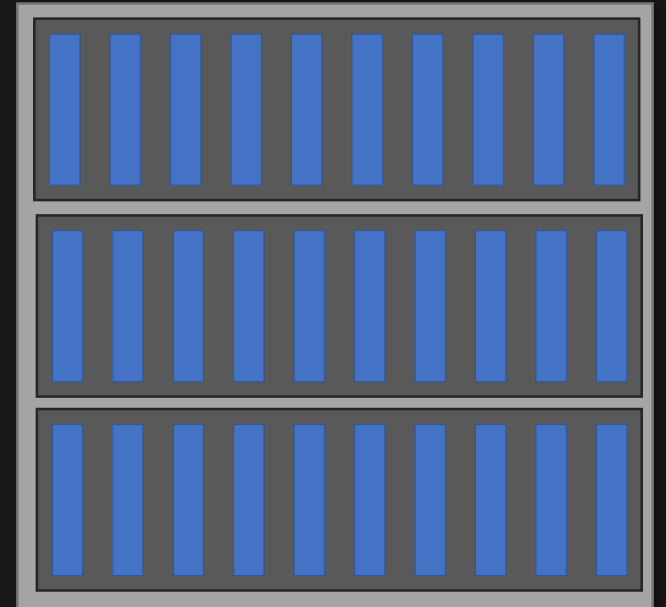

5 – Watch the occupancy!

```
__device__ float inv2x2(const float mat[2][2], float
inv[2][2])
{
    float det = det2x2(mat);
    if(det){
        float ud = 1.0 / det;
        inv[0][0] = mat[1][1] * ud;
        inv[1][0] = -mat[1][0] * ud;
        inv[0][1] = -mat[0][1] * ud;
        inv[1][1] = mat[0][0] * ud;
    }
    else{
        printf("Singular matrix 2x2");
        for(int i = 0; i < 2; i++)
            for(int j = 0; j < 2; j++)
                inv[i][j] = 0.0;
    }
    return det;
}
```



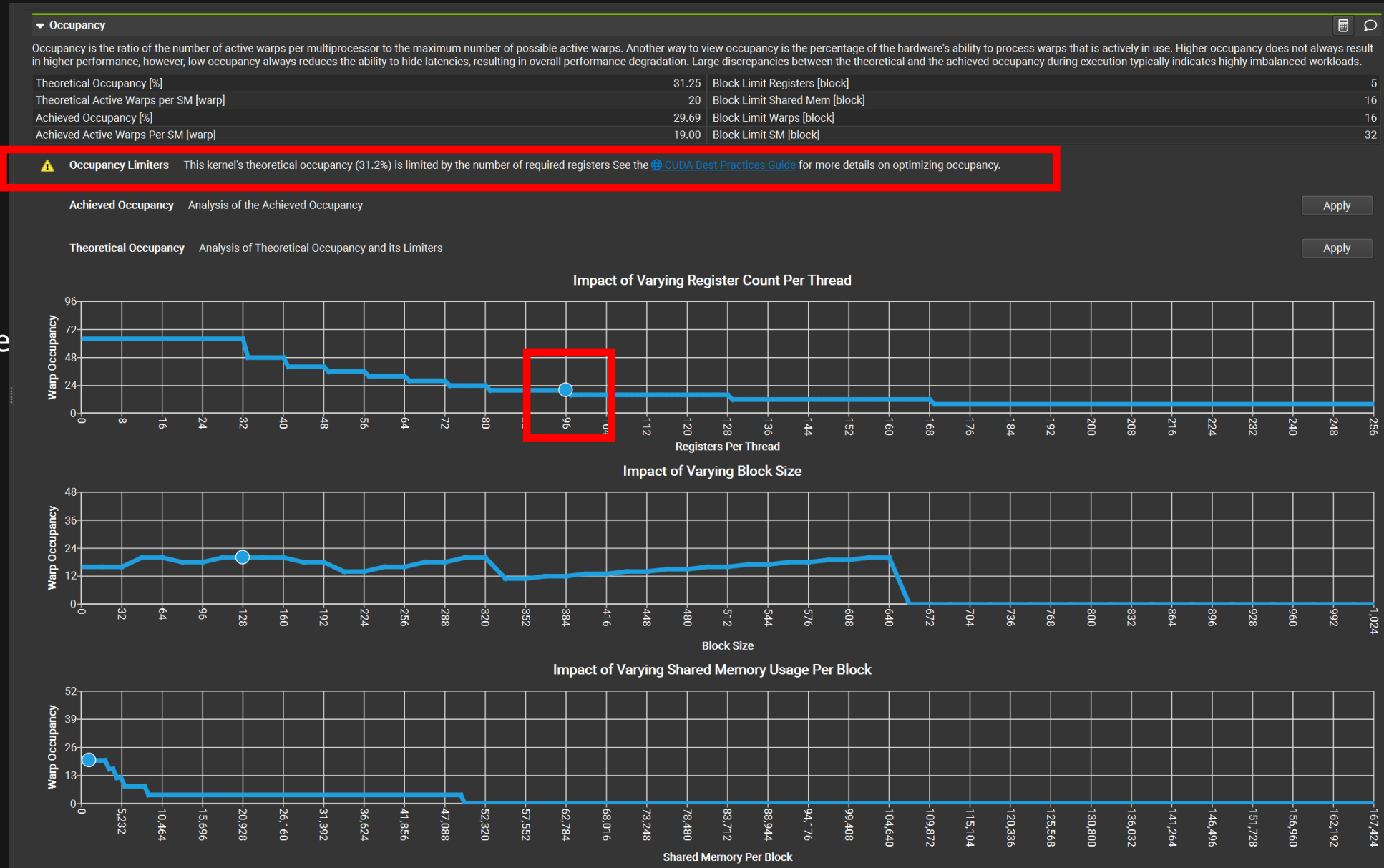
5 – Watch the occupancy!

```
__device__ float inv2x2(const float mat[2][2], float
inv[2][2])
{
    float det = det2x2(mat);
    float ud = 1.0 / det;
    inv[0][0] = mat[1][1] * ud;
    inv[1][0] = -mat[1][0] * ud;
    inv[0][1] = -mat[0][1] * ud;
    inv[1][1] = mat[0][0] * ud;
    return det;
}
```



5 – Is the occupancy a limiting factor?

- ncu will tell you
- Often the case that occupancy is limited by registers, but to be « on a step »



Exercise time!

1. Use floats
2. Use a reordered mesh
3. Transpose the memory access for more coalescence (at the price of locality though!)
4. Catch the remaining double literals
5. Remove the never-accessed debug print code

File : main.h

6 – Kernel fusion

- Kernels comparable in duration
 - Kernels share variables
- Makes sense to merge the kernels to avoid a read and write from global memory



A horizontal timeline diagram consisting of four blue rectangular segments. The first segment is labeled 'axpy', the second is labeled 'dudt', and the third is labeled 'axpy'. The fourth segment is empty. The segments are separated by thin white vertical lines.

axpy

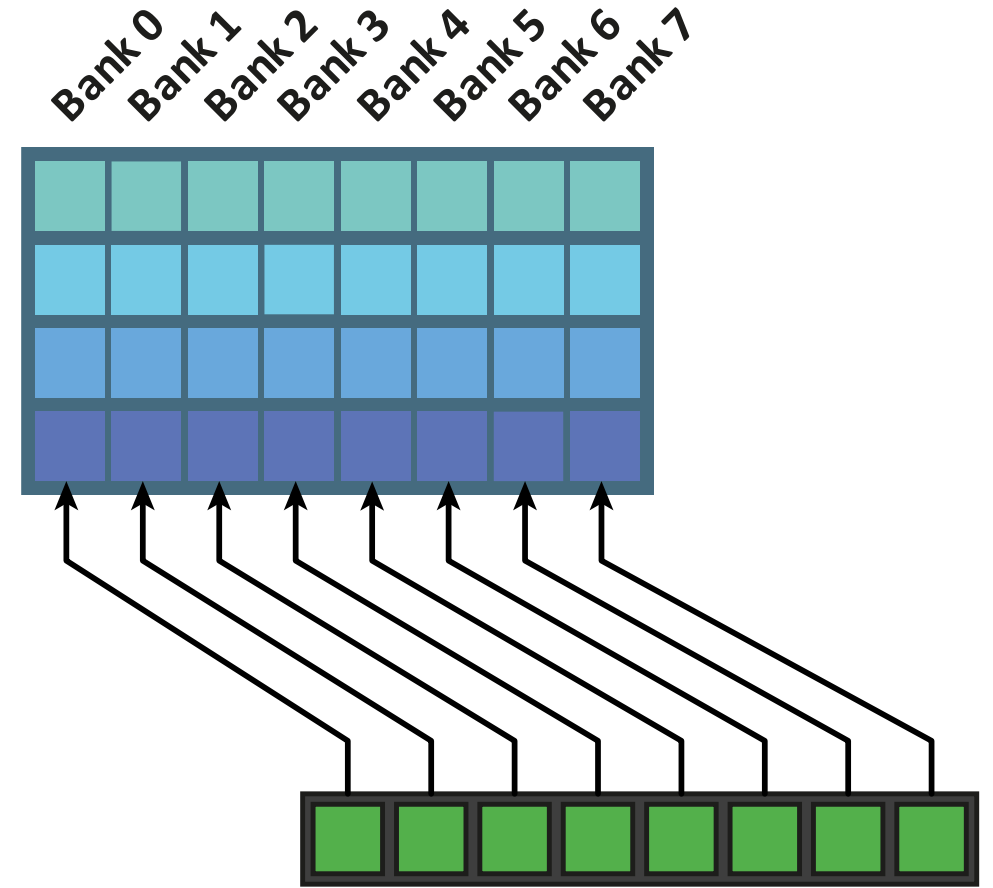
dudt

axpy

Timeline

7 – Using shared memory

- To share data among threads
 - **To manually cache some frequently used data**
 - To reduce register pressure/ local memory usage
 - To allow communication/ data exchange within a group
-
- Shared memory organized in « Banks »
 - Simultaneous accesses to the same bank are serialized
 - Consecutive threads should access consecutive banks



8 – Array of struct of arrays

- Array of struct : perfect locality, bad coalescence



- Struct of array : good coalescence, bad locality



- What if we could combine both?



Array-of-struct-of-array layout : good coalescence, fairly good locality

9 – Free performance*?

- `--use_fast_math` : compiler flag that enables **unsafe** and less accurate but sometimes faster math
- `--extra-device-vectorization`
- `__launch_bounds__()` : Tell the compiler the maximum block size at compile time. Allow optimization that can significantly improve the performance, **or sometimes significantly worsen the performance.**

Usage :

```
__launch_bounds__(BLOCK_SIZE)
__global__ void my_kernel(float a, float* data, int n){...
```

*Sometimes

All of these are just ideas, because

Some edged cases encountered

- `__launch_bounds__()` slowing down the code
- Optimizing for performance is complicated. So complicated that the compiler itself sometimes gets confused
- Adding a big bloc of code made everything 1.2x faster, even when it was never accessed.

Conclusions :

- **profile, benchmark** and run your code before and you try to optimize it. You never know.
- **Use the tools available**, ncu is an excellent tool that will also teach you how to do things right