# Profiling and debugging CUDA

Luca Parisi, EPCC, The University of Edinburgh

l.parisi@epcc.ed.ac.uk

11 Jun 2024

www.archer2.ac.uk

# Outline

- Dubbing CUDA with DDT
  - hardware info
  - analyzing a snapshot of the kernel state
- Profiling CUDA with MAP
  - Analise memory transfers
  - Analise warp executions

DDT

# DDT



Run as usual

Select CUDA runtime

Limited support for AMD GPU

# Compiling

Before running DDT you need to compile with the appropriate flags.

```
nvcc -g -G kernel.cu
```

Debugging information on the device
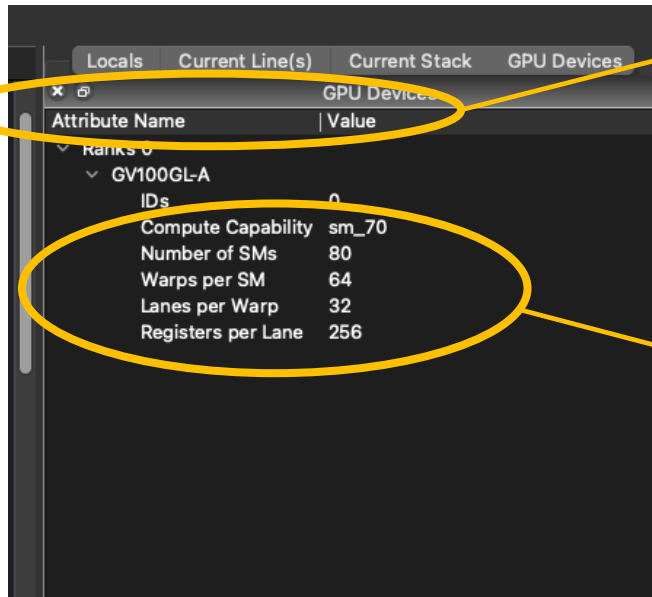
Debugging information on the host

```
nvfortran -O2 -g -gpu=debug kernel.f90
```

Equivalent to `-G` option for NVIDIA C compiler

# DDT

Show information about the device

Which rank is using which GPU ?

Easy access to hardware information needed for analyzing performance.

80 * 64 * 32 = 1.6384 10^5 threads can be on the device

# Stepping trough threads



Advance, break  and pause inside a CUDA kernel

Check number of threads launched.

3.2768 10^6 threads,

x20 the capacity of the device

When execution is paused switch between threads

Threads that have already executed or are yet to be executed

Green: Threads Executing on the device

Blue: Selected Thread

# DDT

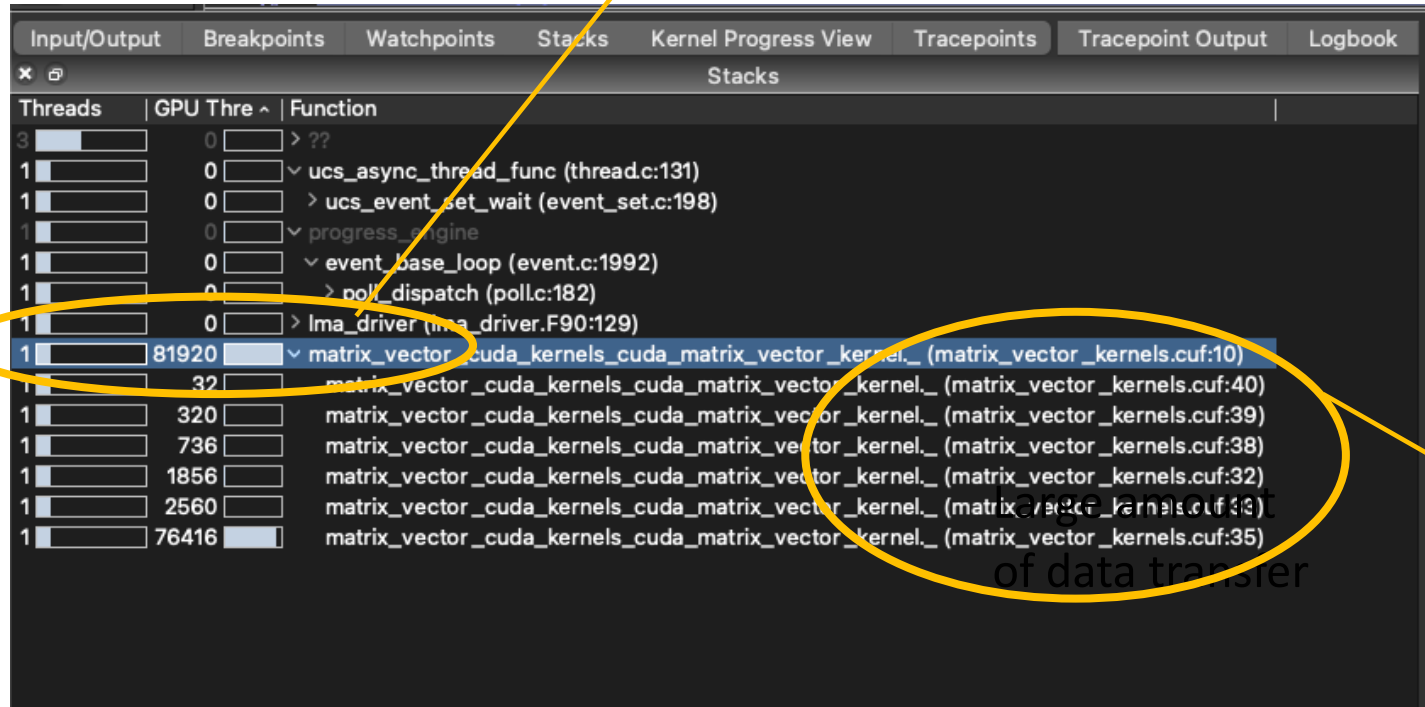81920 threads are executing. About 50% of the device capacity.



Large amount of data transfer

Where are GPU threads ?

MAP

# Launching map

Kernel and transfer analysis each add noticeable overhead

Profile the kernel: warp stall information

```
map --cuda-kernel-analysis --cuda-transfer-analysis -n 1 --mpi=slurm --
mpiargs="--hint=nomultithread --distribution-block:block --cpus-per-task=1" --
profile ./test
```
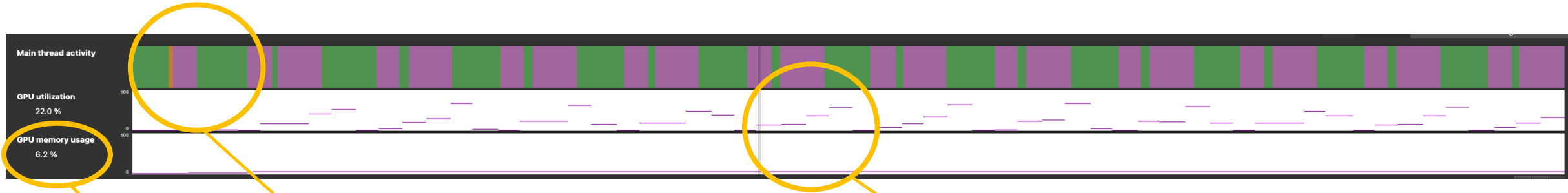
Analyze memory transfer

Memory and Kernel analysis are optional. Can add significant overhead

# Hotspots

Mix of GPU and CPU computation

Green: computation time

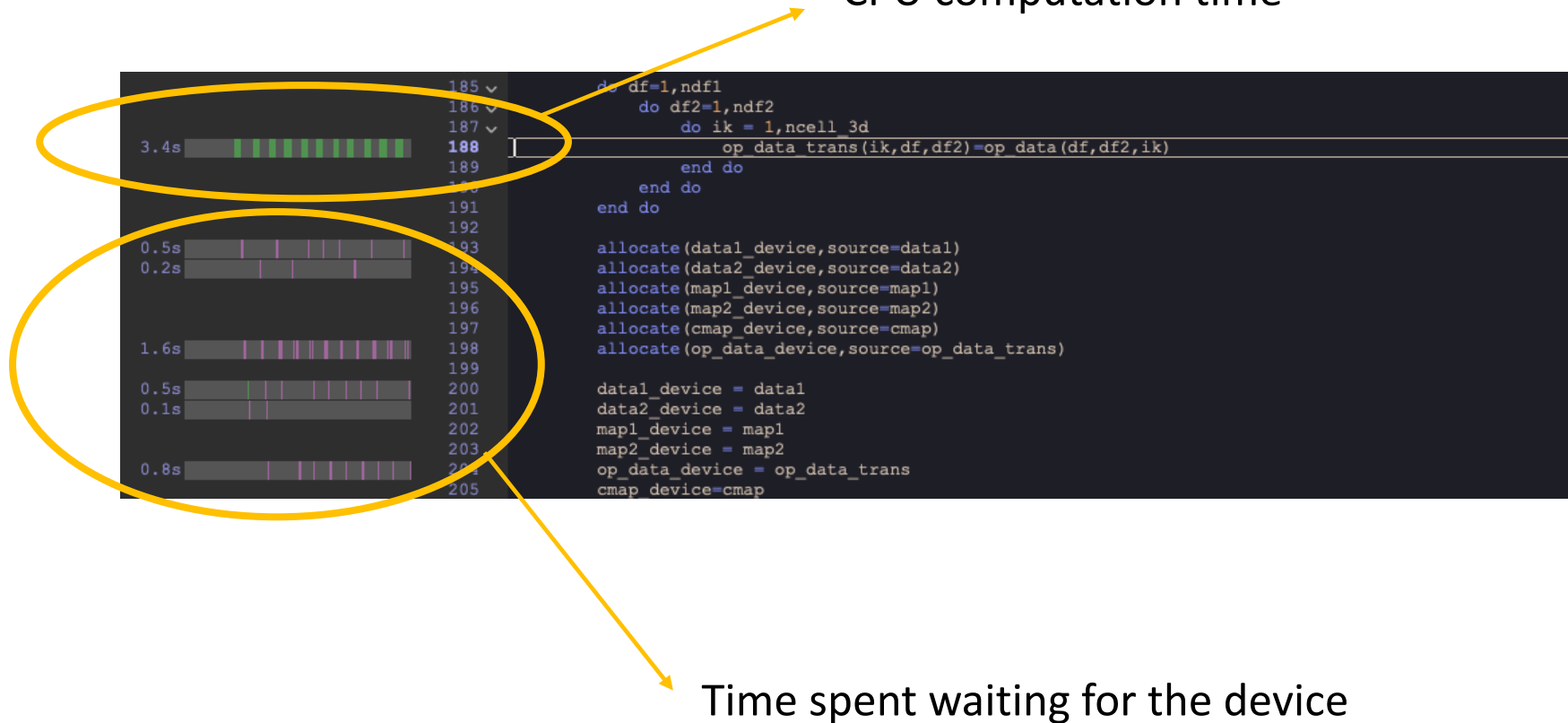Violet: Time spent waiting on the device

Orange: Time spent in IO

Fraction of the time

spent in GPU utilization. Sensitive to undersampling.

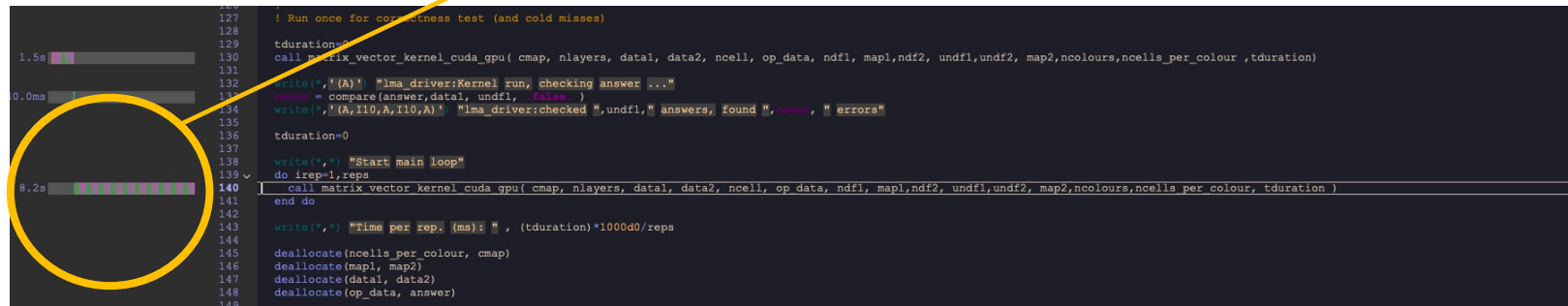Law fraction of the GPU global memory is being used

# CPU and GPU concurrent analysis

|epcc|

CPU computation time

```
                    185 ⌄    do df=1,ndf1
                    186 ⌄        do df2=1,ndf2
                    187 ⌄            do ik = 1,ncell_3d
 3.4s ▮▮|||||||||||  188                op_data_trans(ik,df,df2)=op_data(df,df2,ik)
                    189                end do
                    190            end do
                    191        end do
                    192
 0.5s ▮▮|  ||||||   193    allocate(data1_device,source=data1)
 0.2s ▮▮ |     |    194    allocate(data2_device,source=data2)
                    195    allocate(map1_device,source=map1)
                    196    allocate(map2_device,source=map2)
                    197    allocate(cmap_device,source=cmap)
 1.6s ▮▮| ||||||||| 198    allocate(op_data_device,source=op_data_trans)
                    199
 0.5s ▮▮|||||||||   200    data1_device = data1
 0.1s ▮▮||          201    data2_device = data2
                    202    map1_device = map1
                    203    map2_device = map2
 0.8s ▮▮| |||||||||| 204    op_data_device = op_data_trans
                    205    cmap_device=cmap
```

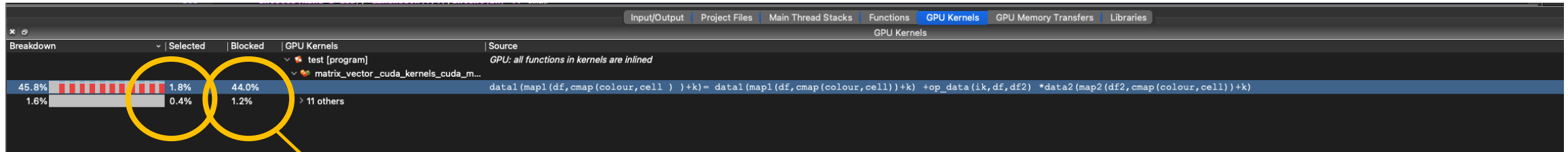Time spent waiting for the device

# Where is the GPU time ?

Mix of GPU and CPU computation



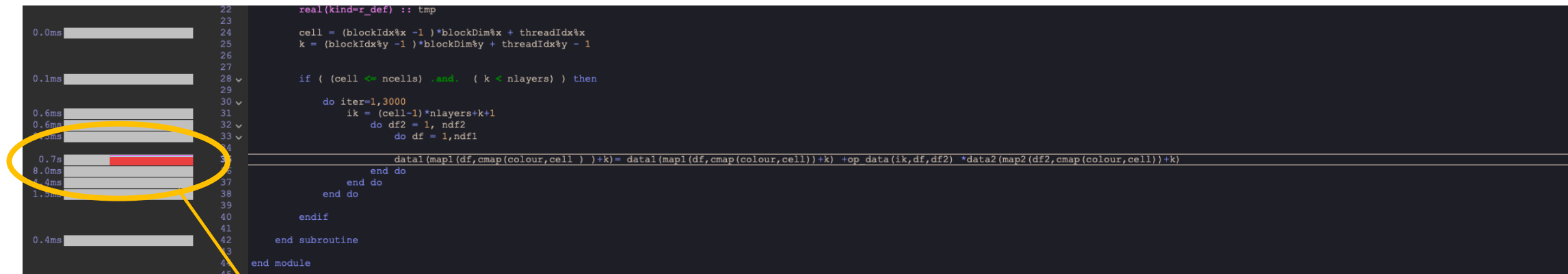Transfers from device after kernel execution
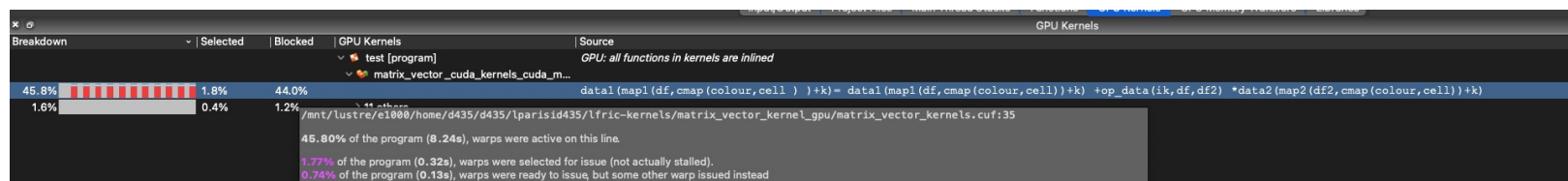
# Kernel execution



Fraction of warps executing
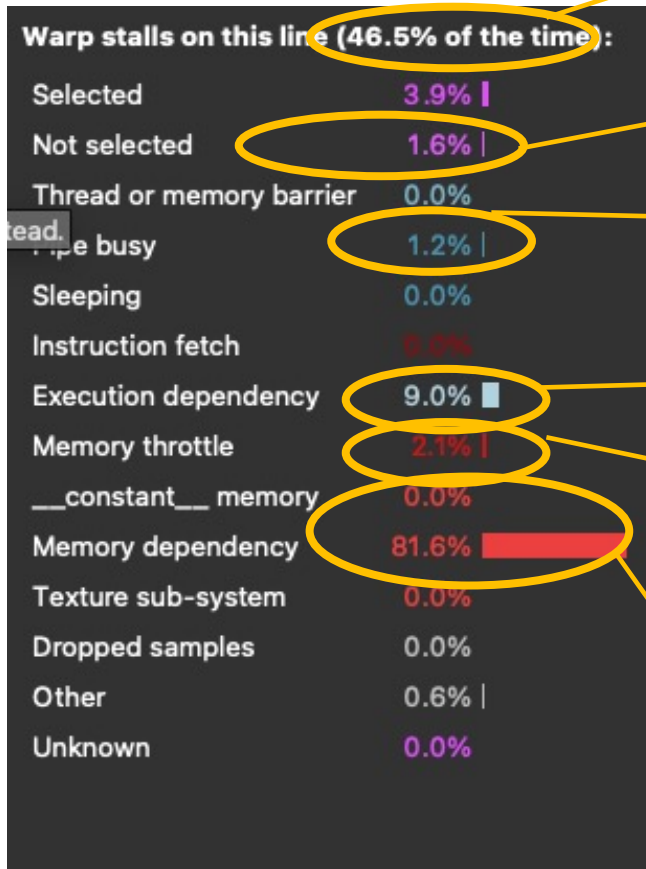
Fraction of warps stalled for some reason

# Kernel execution



The gpu kernel spent the whole time in a kernel

# Kernel execution

| Warp stalls on this line (46.5% of the time): | |
|---|---|
| Selected | 3.9% |
| Not selected | 1.6% |
| Thread or memory barrier | 0.0% |
| Pipe busy | 1.2% |
| Sleeping | 0.0% |
| Instruction fetch | 0.0% |
| Execution dependency | 9.0% |
| Memory throttle | 2.1% |
| __constant__ memory | 0.0% |
| Memory dependency | 81.6% |
| Texture sub-system | 0.0% |
| Dropped samples | 0.0% |
| Other | 0.6% |
| Unknown | 0.0% |

Warps ready for computation, but not executing as not enough resources available

Being executed, but stalling as functional units are not available

Waiting for the result of an arithmetic operation

Too many memory requests

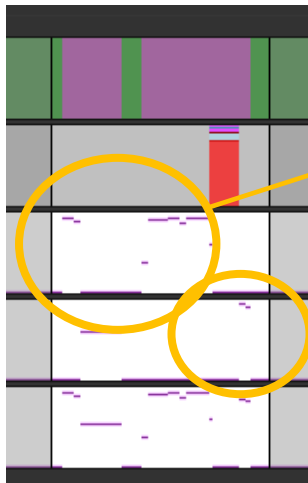Most of the warps waiting for data to be loaded from memory

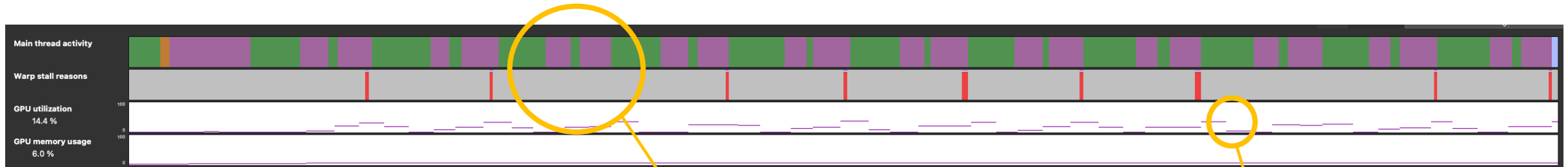# Memory transfers

Many repeated transfers

Transfers to device before kernel execution

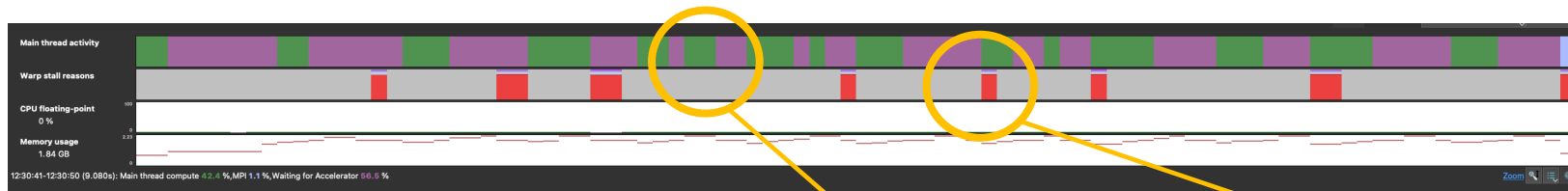Transfers from device after kernel execution

# Sampling



Missing information
about warp stall reasons

GPU usage in
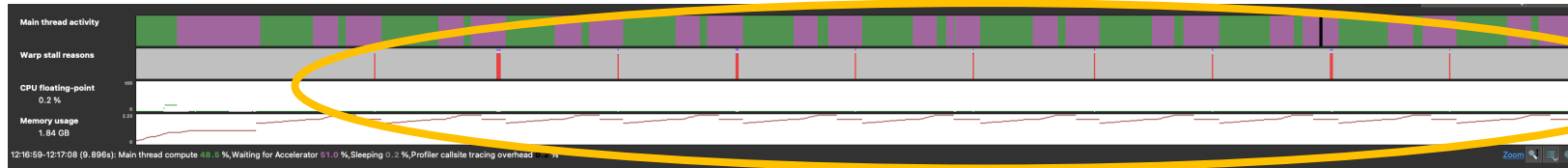unexpected areas of the
code

# Sampling



Kernels in wrong position

Kernel is missing. Time intervals are fluctuating wildly

Inconsistencies and blocky appearance suggests we are not collecting enough samples

# Sampling

```
export FORGE_SAMPLER_INTERVAL=1
export FORGE_SAMPLER_GPU_INTERVAL=1
```
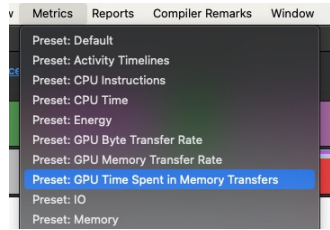


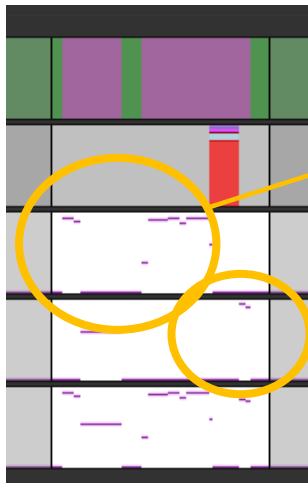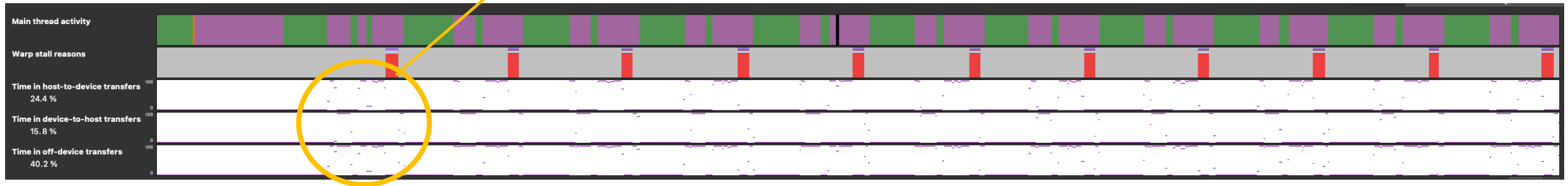Set sampling interval to 1ms, the maximum supported value

Kernels all presents and in the expected position.

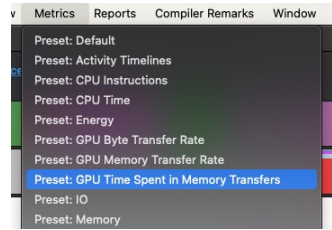Balance between amount of information and sampler overhead
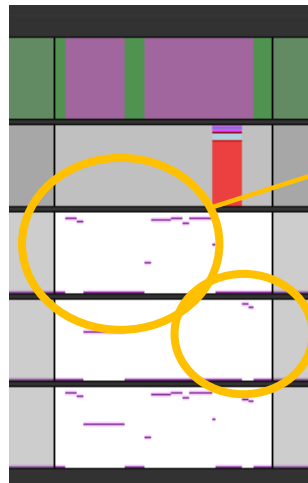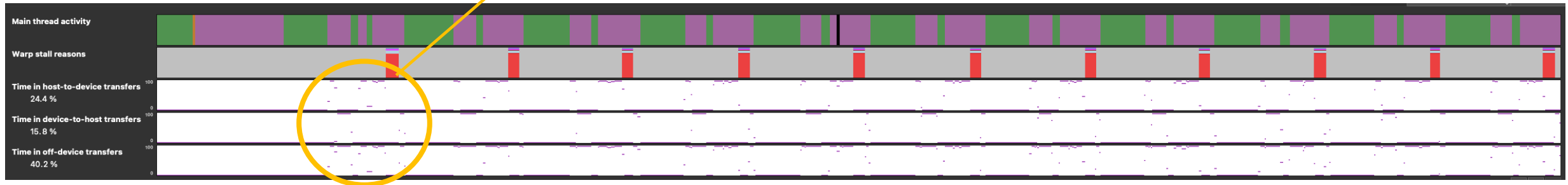
# Memory transfers

Many repeated transfers

Transfers to device before kernel execution

Transfers from device after kernel execution
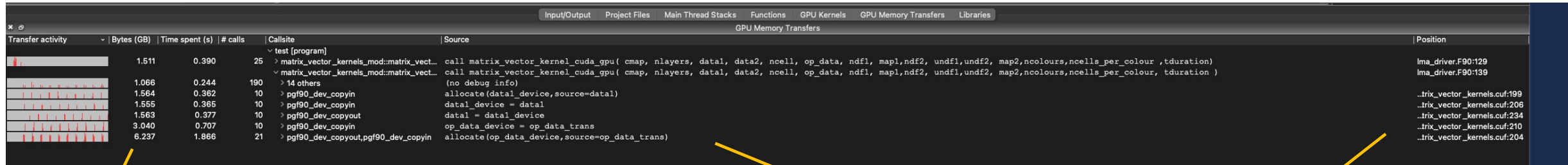
# Memory transfers



Many repeated transfers

Transfers to device before kernel execution

Transfers from device after kernel execution

# Memory transfers



Large amount of data transfer

Sizable amount of time spent in GPU transfers

Source code responsible for the transfer

# Conclusions

- Debug both CPU and GPU with DDT
    - Step into both GPU kernels and CPU kernels
- Profile both CPU and GPU with MAP
    - Time spent waiting on device
    - Memory transfers
    - Stalling reasons