

Miniweather Application GPU Acceleration

Michel Herrera Sanchez

DLI Instructor Evaluation OpenACC

GPU Computing

Differences between CPU and GPU Architecture

CPU Architecture

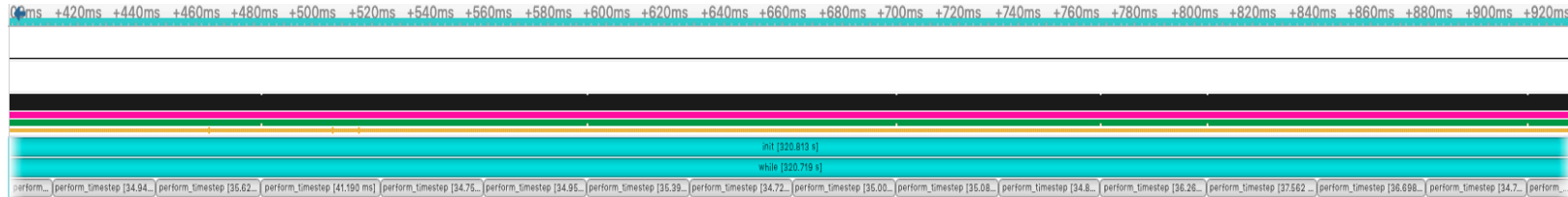
- Designed for Latency workloads
- Lower density of cores concentrations inside the CPU die
- Faster processors speed
- Large Main Memory Capacity and larges L2/L1 Caches
- Complex hardware with Branch Predictor
- Data Parallel Execution using SIMD AVX-512

NVIDIA GPU Architecture

- Designed for high throughput workloads
- High density of cores concentrations inside the GPU die (thousands)
- Slower processors speed
- High Bandwidth Memory Access
- GPU Memory is not integrated in the Host and requires data transfer
- Data Parallel Execution using SIMT Warp Size 32

GPU Computing

Expected speedup



By Amdahl Law : $\text{speedup} = 1 / (s + p / N)$

s -> serial proportion = $320.813 / 641.532 = 0.50$

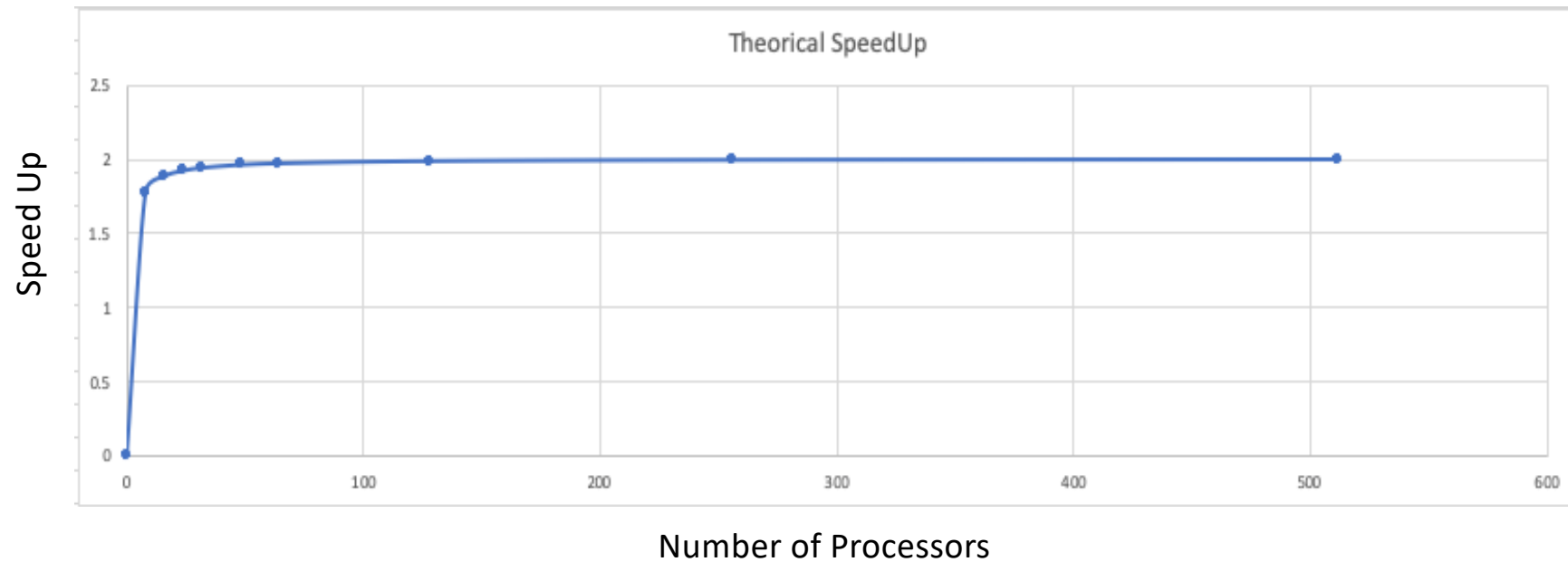
p -> parallelized part = $1 - s = 0.50$

N -> number of processors

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Range
33.4%	320.813 s	1	320.813 s	320.813 s	320.813 s	320.813 s	320.813 s	0 ns init
33.3%	320.719 s	1	320.719 s	320.719 s	320.719 s	320.719 s	320.719 s	0 ns while
33.3%	320.280 s	9000	35.587 ms	35.124 ms	33.589 ms	44.890 ms	1.361 ms	perform_timestep

GPU Computing

Expected speedup



The expected speed up of MiniWeather application is limited to 2x no matter how many processors (N) are used to run the parallelized proportion.

GPU Computing

Is the problem Memory Bound or Compute Bound?

The compute-to-global-memory-access-ratio allow us determine if the application is memory bound or compute bound

compute-to-global-memory-access-ratio = # floating point operations / memory access

The dominant kernel in the Miniweather application parallelized fraction is :

- **compute_tendencies_x**

```
//Fourth-order-accurate interpolation of the state
vals[ll] = -stencil[0] / 12 + 7 * stencil[1] / 12 + 7 * stencil[2] / 12 - stencil[3] / 12;
//First-order-accurate interpolation of the third spatial derivative of the state (for artificial viscosity)
d3_vals[ll] = -stencil[0] + 3 * stencil[1] - 3 * stencil[2] + stencil[3];
```

stencil[] is a local array variable but its value is first loaded from memory before usage

compute-to-global-memory-access-ratio for **compute_tendencies_x** kernel is $10 / (4.0 * 2) \rightarrow 1.25$

GPU Computing (V100)

Is the problem Memory Bound or Compute Bound?

Miniweather **compute_tendencies_x**

Compute-Global-Memory-Access-Ratio: 1.25

Applications uses Double precision

NVIDIA GPU V100

Global Memory Bandwidth: 900 GB/s

Double Precision: 7 TFLOPS

Peak Compute-To-Global-Memory-Access-Ratio : 7.7

With a Compute-To-Global-Memory-Access-Ratio of 1.25 the kernel **compute_tendencies_x** throughput will be limited on how fast operands are fetched from memory. This make the dominant kernel **compute_tendencies_x** /application memory-bound. In fact the kernel when ported to GPU V100 will not achieve more than 112.5 GFLOPS.

Multicore Parallelization Results

Application	target	hardware	Time (s)	Speedup
Serial		intel	307.842808	
OpenACC	multicore	Intel 8 cores	82.222548	3.74

Using default values

nx_glob, nz_glob: 400 200 dx,dz: 50.000000 50.000000

dt: 0.166667

Flags to Compile the Code

```
# Copyright (c) 2020 NVIDIA Corporation. All rights reserved.
CC := pgc++
CFLAGS := -O3 -w
# ACCFLAGS := -ta=tesla -Minfo=accel
# ACCFLAGS := -ta=tesla:managed -Minfo=accel
ACCFLAGS := -ta=multicore -Minfo=accel
LDFLAGS := -lnetcdf -ldl
NVTXLIB := -I/opt/nvidia/hpc_sdk/Linux_x86_64/20.9/cuda/11.0/include

miniWeather: miniWeather_openacc.cpp
——»${CC} ${CFLAGS} ${ACCFLAGS} -o miniWeather miniWeather_multicore.cpp ${NVTXLIB} ${LDFLAGS}

clean:
——»rm -f *.o miniWeather
```


OpenACC directives / optimization

Application Code Sample

```
420 #pragma acc parallel loop private(indt, indf1, indf2, inds)
421 for (ll = 0; ll < NUM_VARS; ll++)
422 {
423     #pragma acc loop vector collapse(2)
424     for (k = 0; k < nz; k++)
425     {
426         for (i = 0; i < nx; i++)
427         {
428             indt = ll * nz * nx + k * nx + i;
429             indf1 = ll * (nz + 1) * (nx + 1) + (k) * (nx + 1) + i;
430             indf2 = ll * (nz + 1) * (nx + 1) + (k + 1) * (nx + 1) + i;
431             tend[indt] = -(flux[indf2] - flux[indf1]) / dz;
432             if (ll == ID_WMOM)
433             {
434                 inds = ID_DENS * (nz + 2 * hs) * (nx + 2 * hs) + (k + hs) * (nx + 2 * hs) + i + hs;
435                 tend[indt] = tend[indt] - state[inds] * grav;
436             }
437         }
438     }
439 }
```

OpenACC directives / optimization

Application Code Sample Description

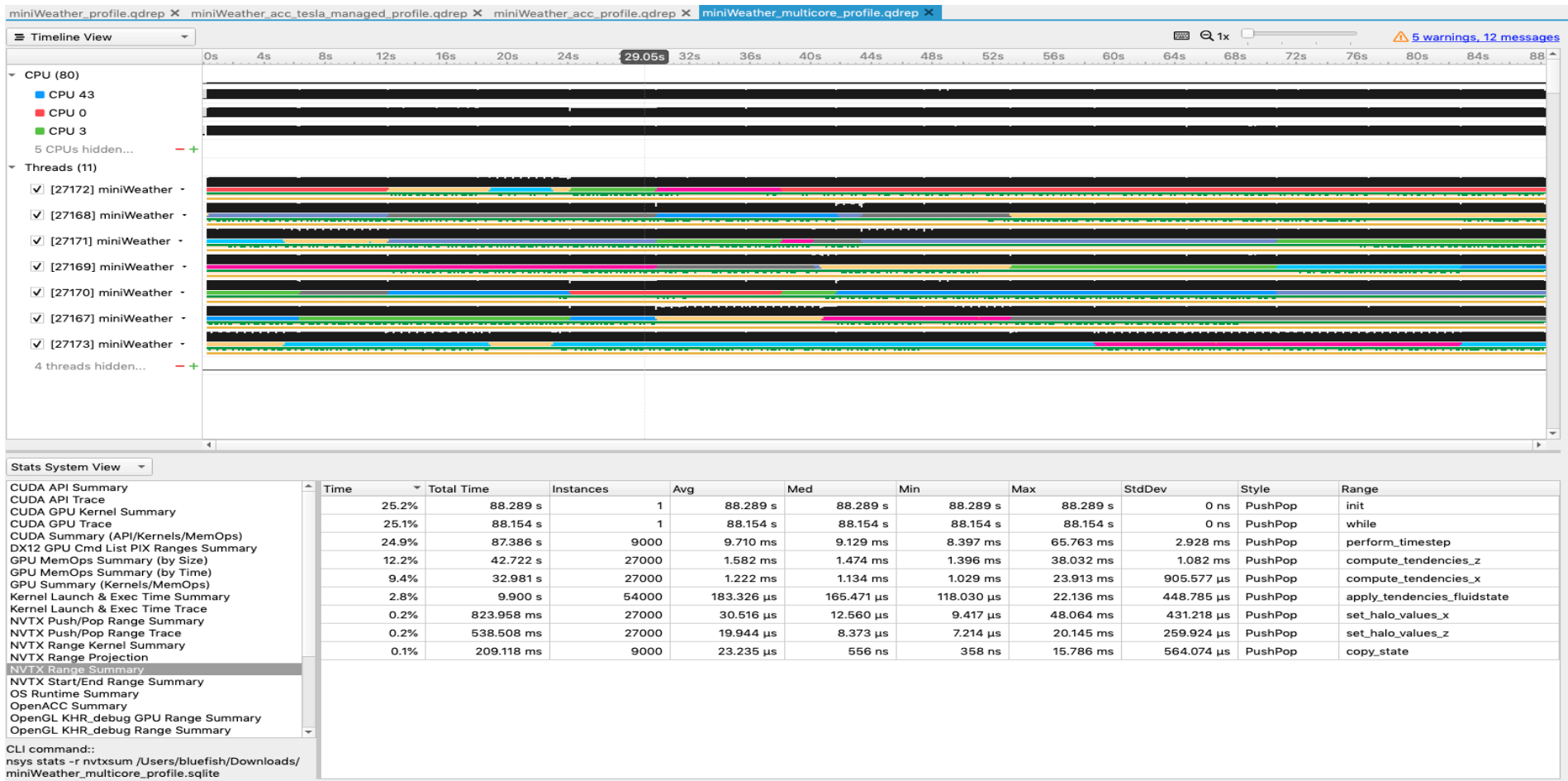
External loop in line 421 is parallelized and distributed among a gang of threads. This is done using the openacc directive `#pragma acc parallel loop`.

Variables that could be potentially shared by multiples threads are declared private and local to each thread in the gang. Line 420

Inner loops at line 424 and 426 are collapsed in one loop and vectorized.

OpenACC directives / optimization

Multicore nsys system profiler output and application timeline



GPU Parallelization Results

Application	target	hardware	Time (s)	Speedup
Serial		intel	307.842808	
OpenACC	multicore	Intel 8 cores	82.222548	3.74
OpenACC	tesla:managed	GPU V100	11.630973	26.46
OpenACC	tesla	GPU V100	12.010639	25.63

Using default values

nx_glob, nz_glob: 400 200 dx,dz: 50.000000 50.000000

dt: 0.166667

Flags to Compile the Code

```
# Copyright (c) 2020 NVIDIA Corporation. All rights reserved.
CC := pgc++
CFLAGS := -O3 -w
ACCFLAGS := -ta=tesla -Minfo=accel
# ACCFLAGS := -ta=tesla:managed -Minfo=accel
# ACCFLAGS := -ta=multicore -Minfo=accel
LDFLAGS := -lnetcdf -ldl
NVTXLIB := -I/opt/nvidia/hpc_sdk/Linux_x86_64/20.9/cuda/11.0/include

miniWeather: miniWeather_openacc.cpp
——»${CC} ${CFLAGS} ${ACCFLAGS} -o miniWeather miniWeather_openacc.cpp ${NVTXLIB} ${LDFLAGS}

clean:
——»rm -f *.o miniWeather
```

OpenACC directives / optimization

Data Optimization directives to minimize data transfer between host and device when not using managed memory. CUDA kernels don't need to wait for data transfer during their execution. Present clause using must of the time.

Parallel Loop optimizations and vector size adjusted to 32.

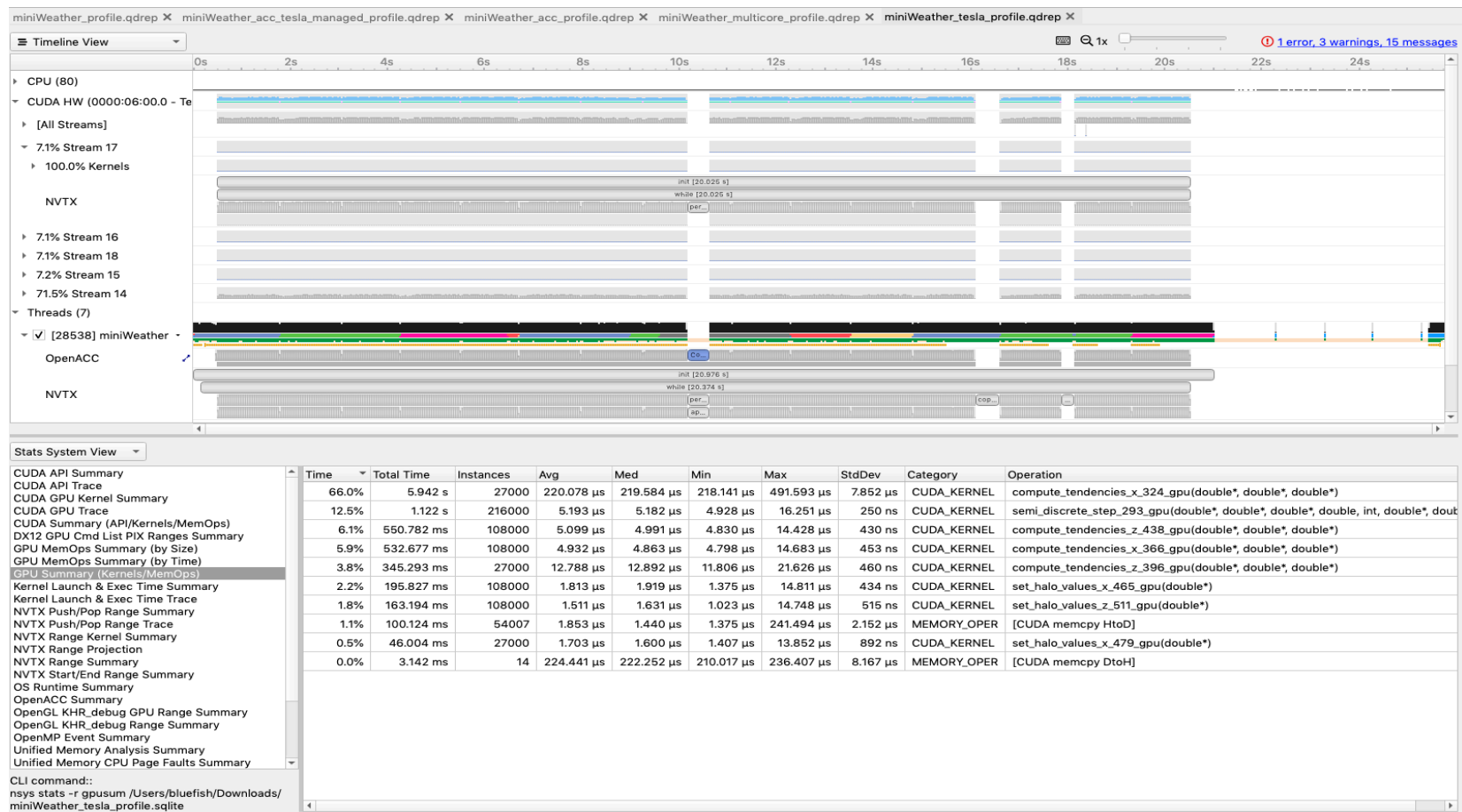
Variables that could potentially be shared among multiples threads are declared private and local to each thread in the gang.

Inner loops are collapsed in one loop and vectorized.

Independent Loop iterations are parallelized using asynchronous cuda streams.

OpenACC directives / optimization

Application Timeline and GPU summary (Kernels/Memory).



OpenACC directives / optimization

Private variables definition per thread to avoid race conditions

Data already on device (present)

```
287 //Apply the tendencies to the fluid state
288 nvtxRangePushA("apply_tendencies_fluidstate");
289
290 for (ll = 0; ll < NUM_VARS; ll++)
291 {
292     #pragma acc parallel loop collapse(2) private(indt, inds) present( state_init[0:(nx + 2 * hs) * (nz + 2 * hs) * NUM_VARS],
state_out[0:(nx + 2 * hs) * (nz + 2 * hs) * NUM_VARS], tend[0:nx * nz * NUM_VARS ]) vector_length(32) async(ll)
293     for (k = 0; k < nz; k++)
294     {
295         for (i = 0; i < nx; i++)
296         {
297             inds = ll * (nz + 2 * hs) * (nx + 2 * hs) + (k + hs) * (nx + 2 * hs) + i + hs;
298             indt = ll * nz * nx + k * nx + i;
299             state_out[inds] = state_init[inds] + dt * tend[indt];
300         }
301     }
302 }
303 #pragma acc wait
304 nvtxRangePop();
305 }
306 }
```

Independent loop iterations are parallelized using asynchronous CUDA streams.

Parallel loop optimization collapse and vectorization

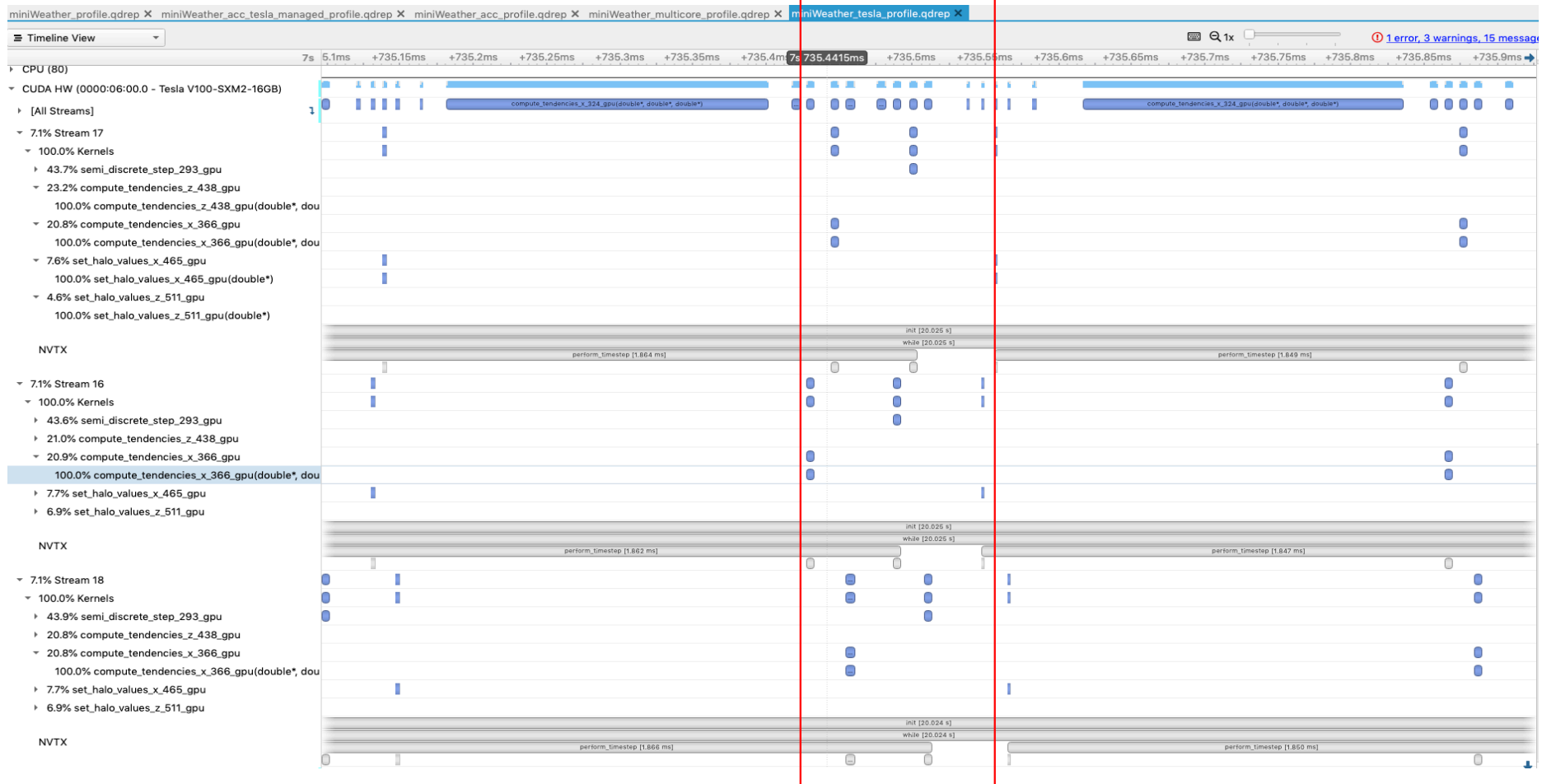
OpenACC directives / optimization

```
In [9]: !cd ../source_code && make clean && make
```

```
rm -f *.o miniWeather
pgc++ -O3 -w -ta=tesla -Minfo=accel -o miniWeather miniWeather_openacc.cpp -I/opt/nvidia/hpc_sdk/Linux_x86_64/20.9/cuda/11.0/include -lnetcdf -ldl
main:
    167, Generating copyin(hy_pressure_int[:nz+1],state[:((nz+4)*(nx+4))*4],tend[:(nz*nx)*4],state_tmp[:((nz+4)*(nx+4))*4],hy_dens_cell[:nz+4],hy_dens_theta_cell[:nz+4],flux[:((nz+1)*(nx+1))*4]) [if not already present]
    196, Generating update self(state[:((nz+4)*(nx+4))*4])
semi_discrete_step(double *, double *, double *, double, int, double *, double *):
    293, Generating present(state_init[:((nz+4)*(nx+4))*4],tend[:(nz*nx)*4],state_out[:((nz+4)*(nx+4))*4])
        Generating Tesla code
    293, #pragma acc loop gang, vector(32) collapse(2) /* blockIdx.x threadIdx.x */
    295, /* blockIdx.x threadIdx.x collapsed */
compute_tendencies_x(double *, double *, double *):
```

Example of Compiler messages.

OpenACC directives / optimization




Kernel `compute_tendencies_x_366_gpu` is launched across multiples streams

OpenACC directives / optimization

Data Optimization directives to minimize data transfer between host and device when not using managed memory. Creation of a data region

```
164     nvtxRangePushA("while");
165
166     #pragma acc data copyin(state[0:(nx + 2 * hs) * (nz + 2 * hs) * NUM_VARS ], state_tmp[0:(nx + 2 * hs) * (nz + 2 * hs) *
NUM_VARS ], flux[0:(nx + 1) * (nz + 1) * NUM_VARS], tend[0:nx * nz * NUM_VARS ], hy_dens_theta_cell[0:nz + 2 * hs],
hy_pressure_int[0:nz + 1], hy_dens_cell[0:nz + 2 * hs])
167     {
168         while (etime < sim_time)
169         {
170             //If the time step leads to exceeding the simulation time, shorten it for the last step
171             if (etime + dt > sim_time)
172             {
173                 dt = sim_time - etime;
174             }
175
176             //Perform a single time step
177             nvtxRangePushA("perform_timestep");
178             perform_timestep(state, state_tmp, flux, tend, dt);
179             nvtxRangePop();
180
181             //Inform the user
182
183             printf("Elapsed Time: %lf / %lf\n", etime, sim_time);
184
185             //Update the elapsed time and output counter
186             etime = etime + dt;
187             output_counter = output_counter + dt;
188             //If it's time for output, reset the counter, and do output
189
190             nvtxRangePushA("copy_state");
191             if (output_counter >= output_freq)
192             {
193                 output_counter = output_counter - output_freq;
194
195                 #pragma acc update self(state[: (nx+2*hs) * (nz+2*hs) * NUM_VARS ])
196                 output(state, etime);
197             }
198             nvtxRangePop();
199         }
200     }
201 }
```

Most of data is copied to the gpu only once



Dominant DtoH memory transfers



OpenACC directives / optimization

Data Optimization directives to minimize data transfer between host and device when not using managed memory.

Time ▾	Total Time	Count	Avg	Med	Min	Max	StdDev	Operation
97.0%	100.124 ms	54007	1.853 μ s	1.440 μ s	1.375 μ s	241.494 μ s	2.152 μ s	[CUDA memcpy HtoD]
3.0%	3.142 ms	14	224.441 μ s	222.252 μ s	210.017 μ s	236.407 μ s	8.167 μ s	[CUDA memcpy DtoH]

Total ▾	Count	Avg	Med	Min	Max	StdDev	Operation
92.75 MiB	54007	1.76 KiB	1.57 KiB	1.57 KiB	2.52 MiB	21.87 KiB	[CUDA memcpy HtoD]
35.21 MiB	14	2.52 MiB	2.52 MiB	2.52 MiB	2.52 MiB	0 B	[CUDA memcpy DtoH]

Most of the memory transfer occurred when data generated in the device was copied to variable state

Code available at:

https://github.com/mahsanchez/mw_openacc/blob/main/miniweather_openacc.cpp

https://github.com/mahsanchez/mw_openacc/blob/main/miniweather_multicore.cpp