

Hands On Session with Nsight Systems and Nsight Compute

By: Brett Neuman & Daniel Howard, Consulting Services Group, CISL & NCAR

bneuman@ucar.edu & dhoward@ucar.edu

Date: June 16th, 2022

In this notebook we explore profiling of the mini-app MiniWeather to present profiling techniques and code examples. We will cover:

- Overview of Profiling and Performance Sampling Tools
 - Typical development workflows with profiling tools
- NSight Compute for Individual GPU Kernel Performance Analysis
 - How to generate ncu reports and command line parameters
 - Overview of GPU kernel profiling data and source code timing heatmaps
 - External resources for interpreting ncu reports data

Head to the [NCAR JupyterHub portal](https://jupyterhub.hpc.ucar.edu/stable) (<https://jupyterhub.hpc.ucar.edu/stable>) and start a JupyterHub session on Casper login (or batch nodes using 1 CPU, no GPUs) and open the notebook in `10_HandsOnNsight/nsys/10_HandsOnNsight_nsys.ipynb`. Be sure to clone (if needed) and update/pull the NCAR GPU_workshop directory.

Use the JupyterHub GitHub GUI on the left panel or the below shell commands
git clone git@github.com:NCAR/GPU_workshop.git
git pull

Workshop Etiquette

- Please mute yourself and turn off video during the session.
- Questions may be submitted in the chat and will be answered when appropriate.
You may also raise your hand, unmute, and ask questions during Q&A at the end of the presentation.
- By participating, you are agreeing to [UCAR's Code of Conduct](https://www.ucar.edu/who-we-are/ethics-integrity/codes-conduct/participants) (<https://www.ucar.edu/who-we-are/ethics-integrity/codes-conduct/participants>).
- Recordings & other material will be archived & shared publicly.
- Feel free to follow up with the GPU workshop team via Slack or submit support requests to support.ucar.edu (<https://support.ucar.edu>)
 - Office Hours: Asynchronous support via [Slack](https://ncargpuusers.slack.com) (<https://ncargpuusers.slack.com>) or schedule a time with an organizer

Notebook Setup

Set the `PROJECT` code to a currently active project, ie `UCIS0004` for the GPU workshop, and `QUEUE` to the appropriate routing queue depending on if during a live workshop session (`gpuworkshop`), during weekday 8am to 5:30pm MT (`gpudev`), or all other times (`casper`). Due to limited shared GPU resources, please use `GPU_TYPE=gp100` during the workshop. Otherwise, set `GPU_TYPE=v100` (required for `gpudev`) for independent work. See [Casper queue documentation \(`https://arc.ucar.edu/knowledge_base/72581396#StartingCasperjobswithPBS-Concurrentresourcelimits`\)](https://arc.ucar.edu/knowledge_base/72581396#StartingCasperjobswithPBS-Concurrentresourcelimits) for more info.

```
In [ ]: export PROJECT=UCIS0004
          export QUEUE=gpudev
          export GPU_TYPE=gp100

          module load nvhpc/22.5 openmpi &> /dev/null
          export PNEDCDF_INC=/glade/u/apps/dav/opt/pnetcdf/1.12.3/openmpi/4.1.4/nvhpc/22.5/include
          export PNEDCDF_LIB=/glade/u/apps/dav/opt/pnetcdf/1.12.3/openmpi/4.1.4/nvhpc/22.5/lib
```

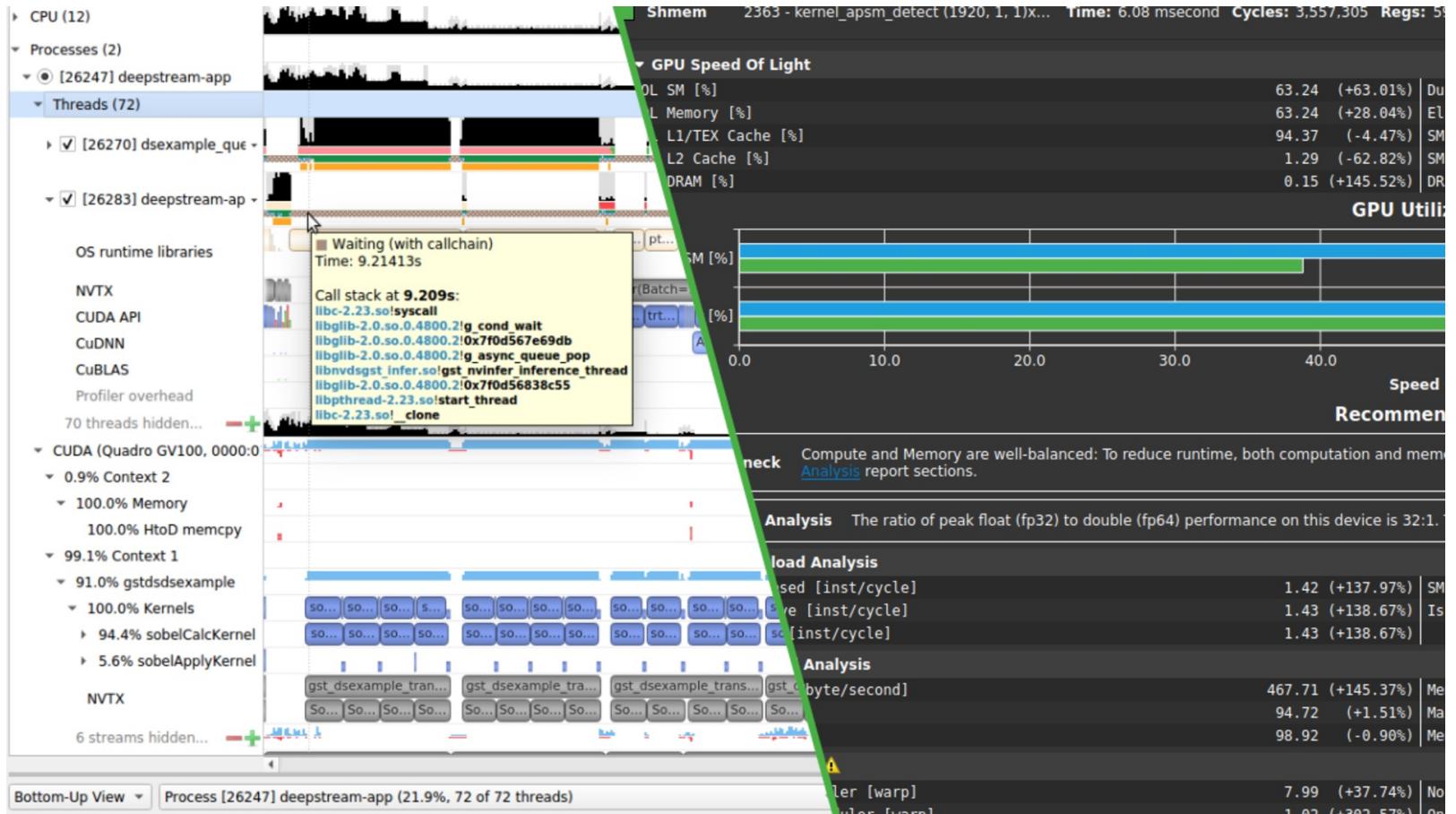
What is a Profiler?

Profilers are tools that **samples** and measure performance characteristics of an executable across its runtime. This information is intended to aid program optimization and performance engineering.

Profiler software that are supported at NCAR include **Arm Map**, **Nsight Systems**, and **Nsight Compute**. All of these tools are able to analyze GPU code. Other profilers you may be aware of include TAU, Intel VTune Advisor, HPC Toolkit, and Vampir.

Today, we will focus on the NVIDIA Nsight profiling tools and usage techniques of these tools.

- **Nsight Systems** - Provides a high level runtime and trace analysis of the program runtime via a measured timeline of various metrics and GPU kernels across a program.
- **Nsight Compute** - Provides an in depth level assessment of individual GPU kernel performance and how various GPU resources are utilized across many different metrics.



Nsight Systems (left) shows a timeline of code runtime.

Nsight Compute (right) records and presents extensive performance statistics for individual kernels.

Profiling Documentation Resources

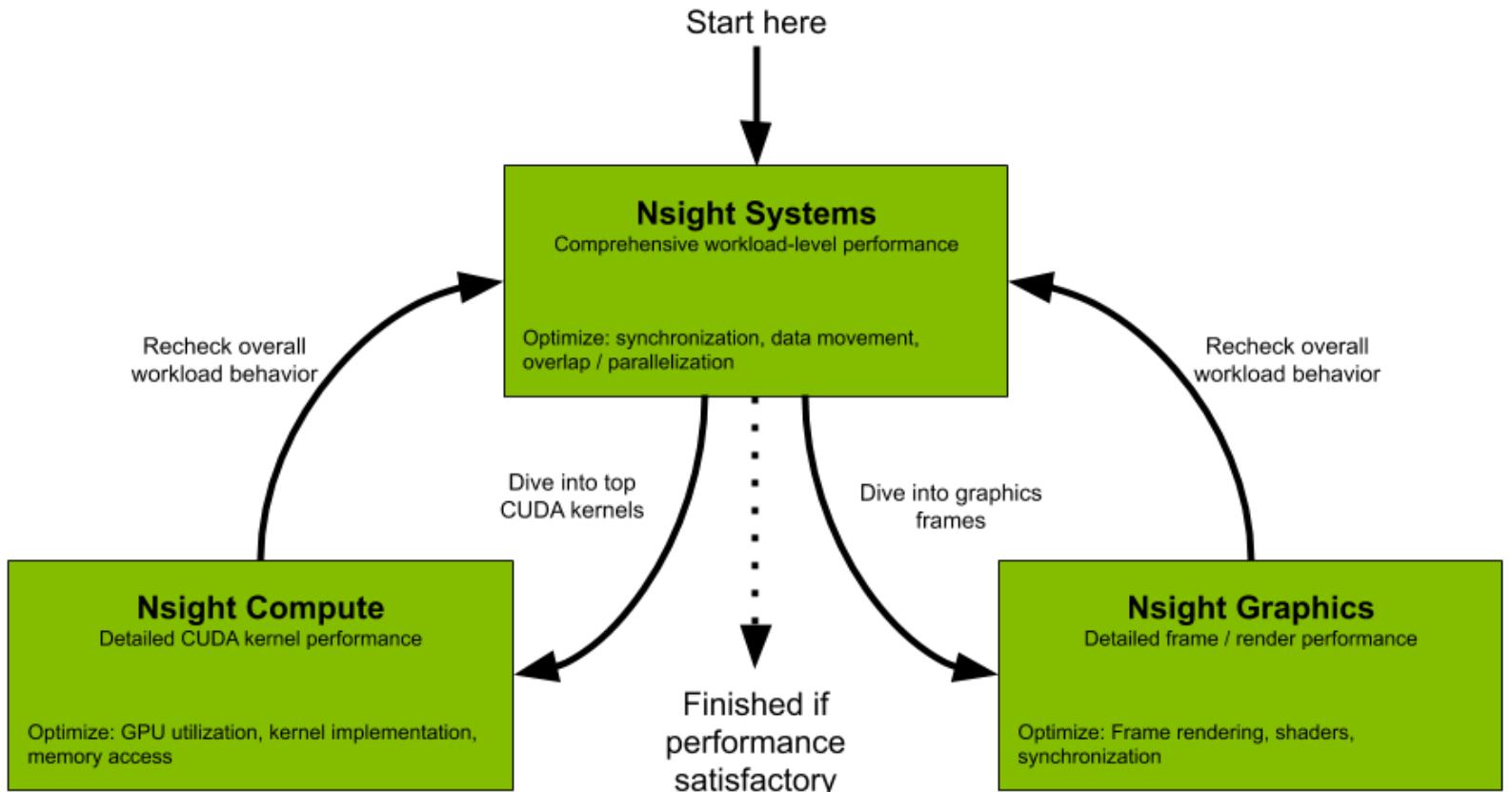
NVIDIA provides extensive documentation for each of these profilers. We will go over basic usage of these tools but to learn more and get the most out of Nsight, consult the below resources:

- [Nsight Systems Main Documentation \(`https://docs.nvidia.com/nsight-systems`\)](https://docs.nvidia.com/nsight-systems)
- [Nsight Compute Main Documentation \(`https://docs.nvidia.com/nsight-compute/`\)](https://docs.nvidia.com/nsight-compute/)
- [Nsight Compute Profiling Guide \(`https://docs.nvidia.com/nsight-compute/ProfilingGuide/index.html`\)](https://docs.nvidia.com/nsight-compute/ProfilingGuide/index.html)
- [Nsight Compute Training Resources \(`https://docs.nvidia.com/nsight-compute/Training/index.html`\)](https://docs.nvidia.com/nsight-compute/Training/index.html) - Forum, Videos, and Blog Posts curated by NVIDIA

An excellent interactive step-by-step tutorial given by Max Katz (NVIDIA) using Nsight Compute to optimize an OpenACC kernel in the BerkeleyGW many-body perturbation theory software can be found at [this Gitlab repository](https://gitlab.com/NERSC/roofline-on-nvidia-gpus) (`https://gitlab.com/NERSC/roofline-on-nvidia-gpus`). A recorded video on this material is [here](https://www.youtube.com/watch?v=fsC3QeZHM1U) (`https://www.youtube.com/watch?v=fsC3QeZHM1U`).

Additionally, the CLI help pages via the `-h` flag for each profiler is a useful quick reference point. Run the below cells to view them.

Profiling Workflow



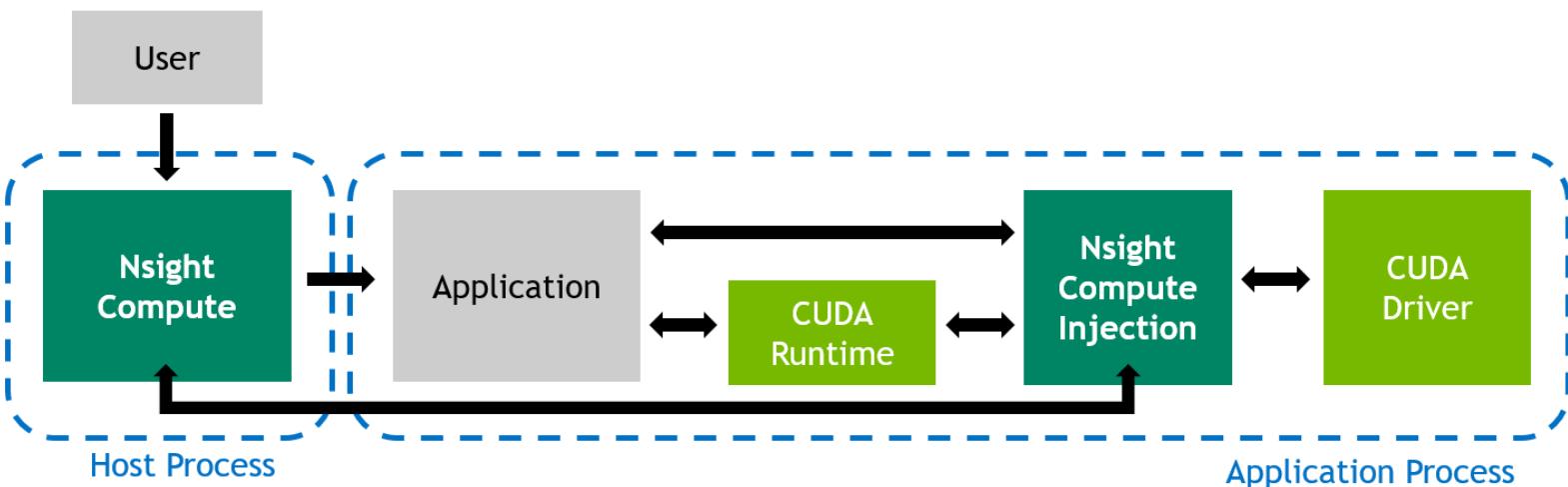
When assessing performance of software, first profile the overall program with Nsight Systems. Then, expensive kernels can be identified and profiled using Nsight Compute.

Iteratively analyze and modify code to optimize performance, up to the amount of effort is worthwhile.

Nsight Compute

After getting a sense of the overall performance of your program with Nsight Systems, use Nsight Compute to dive deeper into the performance of individual GPU kernels.

- CUDA kernel profiler (or CUDA kernels generated by OpenACC/OpenMP/Kokkos code)
- Curates **performance statistics** into targeted metrics sections
- Able to **select amount of data to collect** and how it's presented
 - More detailed analysis has greater **overhead with profiler usage**
- Fully featured **Command Line** and **User Friendly GUI** interfaces
- Regularly updated and customizable Python based rules for **guided analysis** and post-processing



Preparing Code for Nsight Compute

When preparing code for Nsight Compute, an important compile option to add is `-gpu=lineinfo`. **DON'T USE `-pg`, `-g`, or `-G` flags.** The `lineinfo` flag allows the Source/SASS analysis section of Nsight Compute correlate performance information with specific lines of CUDA and/or OpenACC/OpenMP code.

Use the below cell to compile and re-compile MiniWeather after code changes are made. You may also modify the runtime parameters, grid size, and simulation time to investigate how different problem sizes impact performance. Review the generated GPU kernel specifications from the `-Minfo=acc` output.

```
In [ ]: export OPENACC_FLAGS="-acc -gpu=cc60,cc70,lineinfo"

mpif90 -I${PNETCDF_INC} -Mextend -O0 -DN0_INFORM -c miniWeather_mpi_openacc.F90
-o miniWeather_mpi_openacc.F90.o \
-D_NX=9192 -D_NZ=4096 -D_SIM_TIME=0.1 -D_OUT_FREQ=2.0 -D_DATA_SPEC=DATA_SPEC_THERMAL ${OPENACC_FLAGS} -Minfo=acc

mpif90 -Mextend -O3 miniWeather_mpi_openacc.F90.o -o openacc -L${PNETCDF_LIB} -lpnetcdf ${OPENACC_FLAGS}
rm -f miniWeather_mpi_openacc.F90.o
```

Notably, only a short simulation time (enough to cover a few timesteps) is required for us to effectively analyze and optimize model performance.

Nsight Compute CLI Options

- -o <report-name> - Writes output to a *.ncu-rep file to analyze via GUI
 - Without -o, analysis is summarized in stdout.
- -f - Force overwrite of output files
- -c or --launch-count - Specifies the number of kernel launches to profile.
Otherwise, all launched kernels are profiled
- -s or --launch-skip - Skips a specified number of kernel launches. Useful for letting the GPU "warm-up"
- --set <arg> - Sets the amount of data collected and kernel metrics measured, i.e. detailed, full, or others given from --list-sets flag
 - More data collected requires more redundant runs of GPU kernels and increases profiler overhead
- -k or --kernel-name - Specifies the exact name (see nsys) of kernels to be profiled
 - Use -k regex:<expression> to filter kernels by a regex expression
- --nvtx - Enables support for NVTX ranges
- --nvtx-include arg - Filters profiled kernels based on NVTX ranges
- --import-source on - Imports CUDA/source code directly into the report.

Generate Nsight Compute Report

Start with the final version of [miniWeather_mpi_openacc.F90](#) ([miniWeather_mpi_openacc.F90](#)). As we analyze performance, use the generated report to inform code optimizations to experiment with.

First, use the submit script [ncu_bash.sh \(ncu_bash.sh\)](#) to run Nsight Compute on MiniWeather by running command `ncu <ncu options> <exec> <exec arguments>`. Useful ncu options are listed above but also may be reviewed via `ncu -h`.

The first profile run of MiniWeather will profile all kernels using `--sets full` in order to make a baseline (requires redundantly running kernels 73-74 times). When changing code, modify the report filename when you re-run the below cell to help you keep track of reports between different code versions.

In []: `qsub -q $QUEUE -l gpu_type=$GPU_TYPE -A $PROJECT -v NCU_REPORT="MW_DivToMult" ncu_bash.sh`

SHIFT + right click [MW baseline.ncu-report \(MW baseline.ncu-report\)](#) in order to save the Nsight Compute report to your personal machine (or download the file from the left pane explorer). Use your local Nsight Compute client to open the file. Alternatively, after setting `module load nvhpc`, you can run `ncu-ui <report-name>` over a terminal X session or VNC/FastX session on Casper.

Analysis of Nsight Compute Profiles

Depending on the option chosen for `--set` and number of metrics measured, the kernel profiling report will contain a selection of different sections for review covering performance metrics of each kernel profiled.

When using the GUI, **guided analysis** as alerted via exclamation point warning signs will suggest specific issues the profiler identifies and tries to suggest solutions. These are automatically triggered Python rules written by Nsight Compute maintainers and experts, which can be further customized or added to. If you need help interpreting this information, hover your mouse over a piece of information and an informative text box will appear to explain.

Below, we review a few important sections.

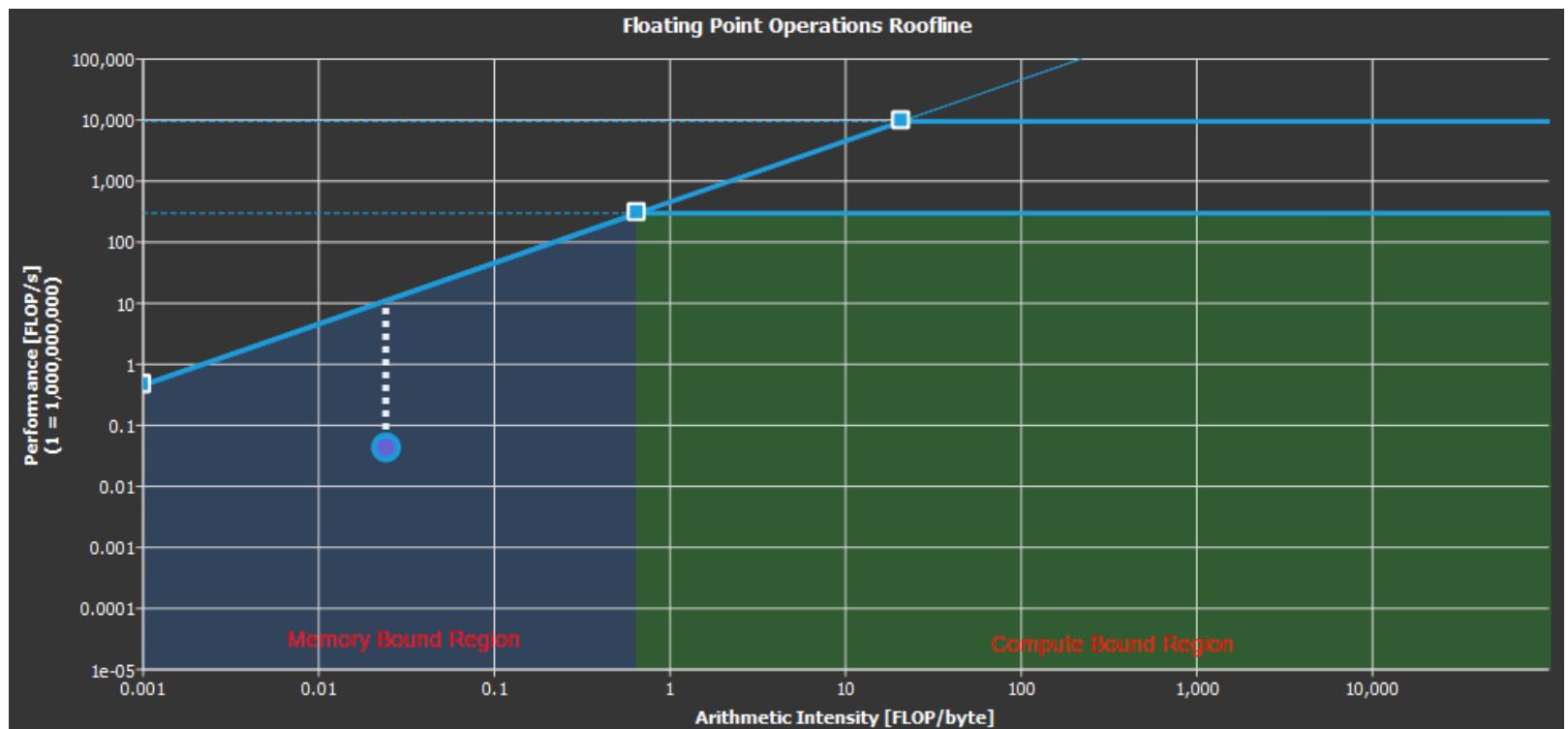
Nsight Compute - GPU Speed of Light



Mouse over each to see the associated metric

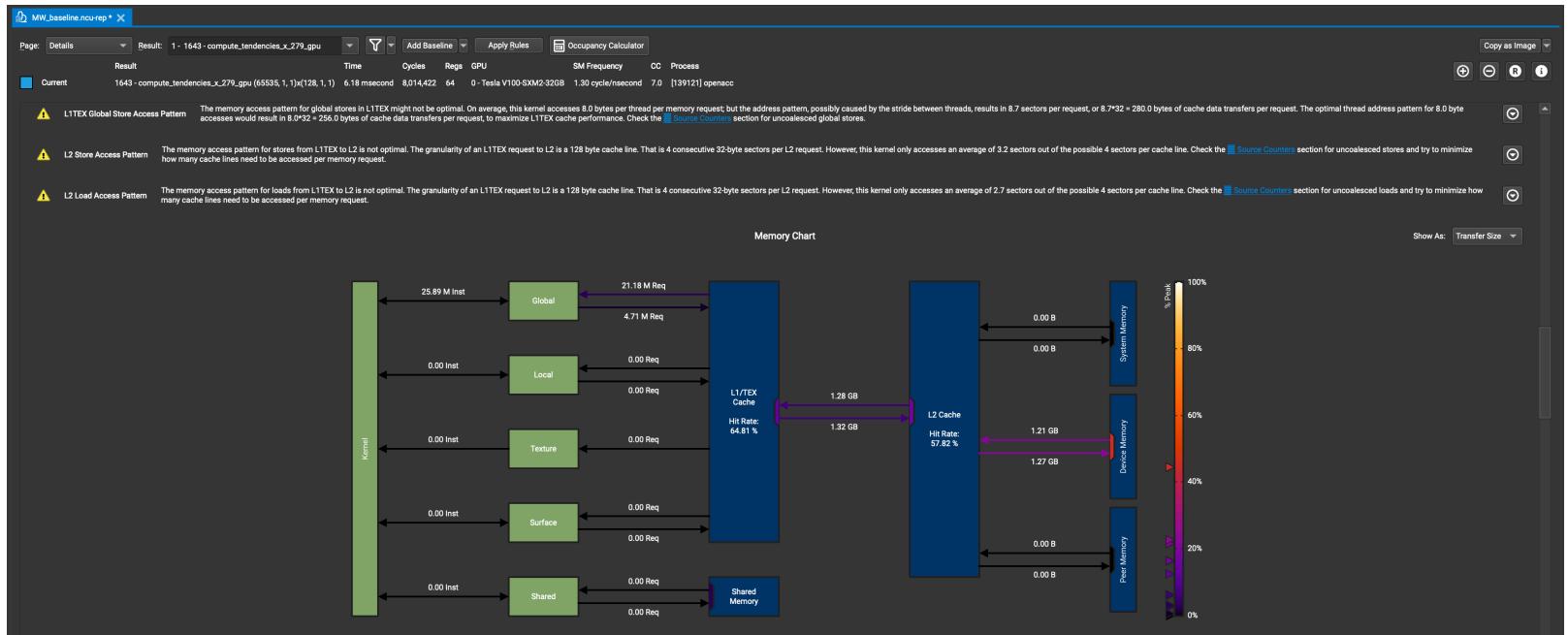
The GPU Speed of Light section highlights to what percentage is this kernel using the full capability of the GPU, both in terms of Streaming Multiprocessor (SM) occupancy and Memory Throughput.

Nsight Compute - Roofline Analysis



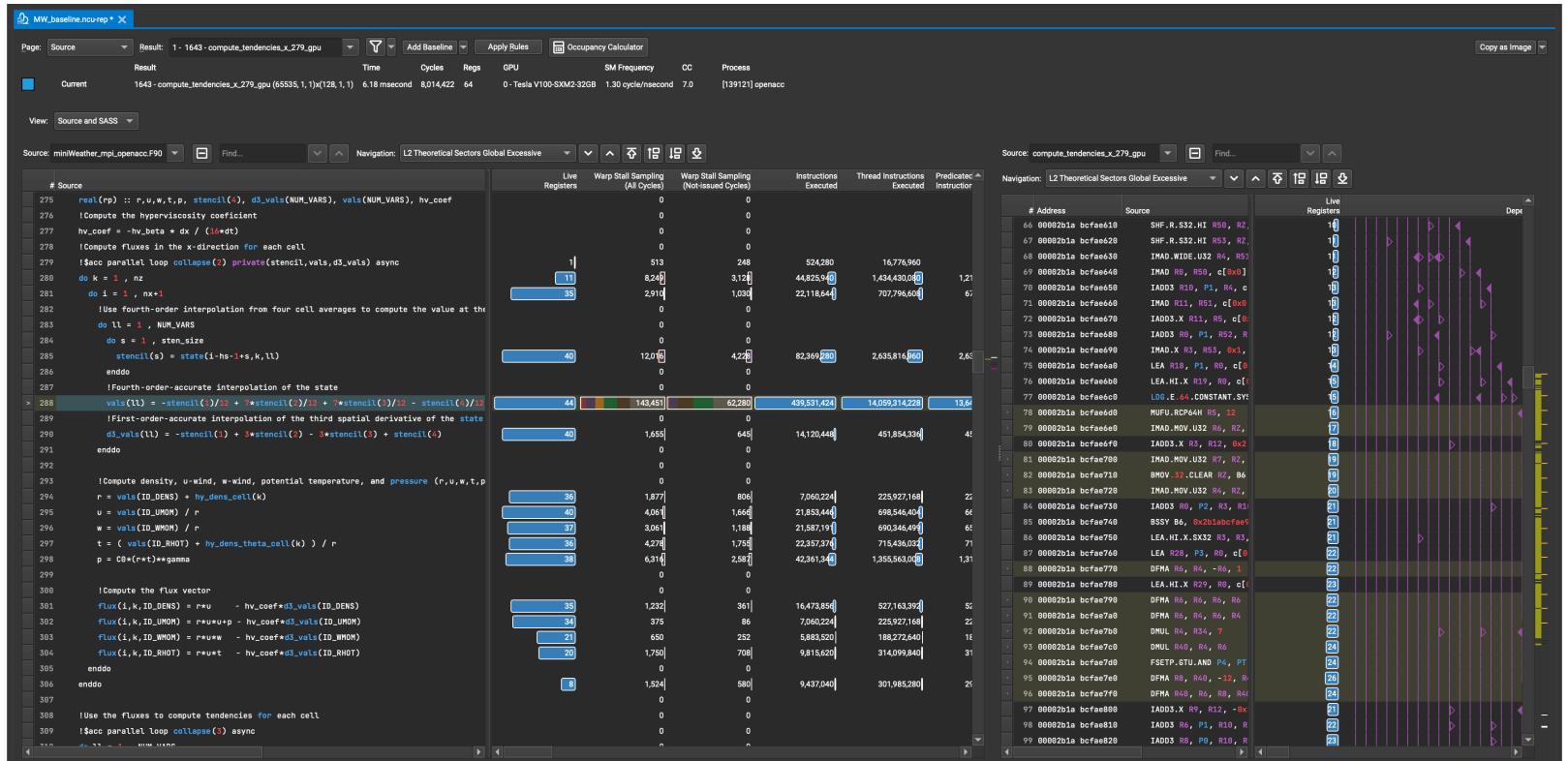
With at least `--set detailed`, a roofline plot is generated, used to determine if the kernel is **compute bound** or **memory bound**. Memory bound kernels can perhaps benefit by assigning more compute operations per thread if possible. Compute bound kernels will likely require further analysis for optimization, typically by checking for warp stalls or coalesced memory issues.

Nsight Compute - Memory Workload Analysis



This section provides a detailed analysis of the memory resources of the GPU. In this case, Nsight Compute identifies that there is an imbalance of data movement between the L1 and L2 caches due to uncoalesced memory. To improve this, memory access patterns need to be re-designed within the source code and OpenACC kernel.

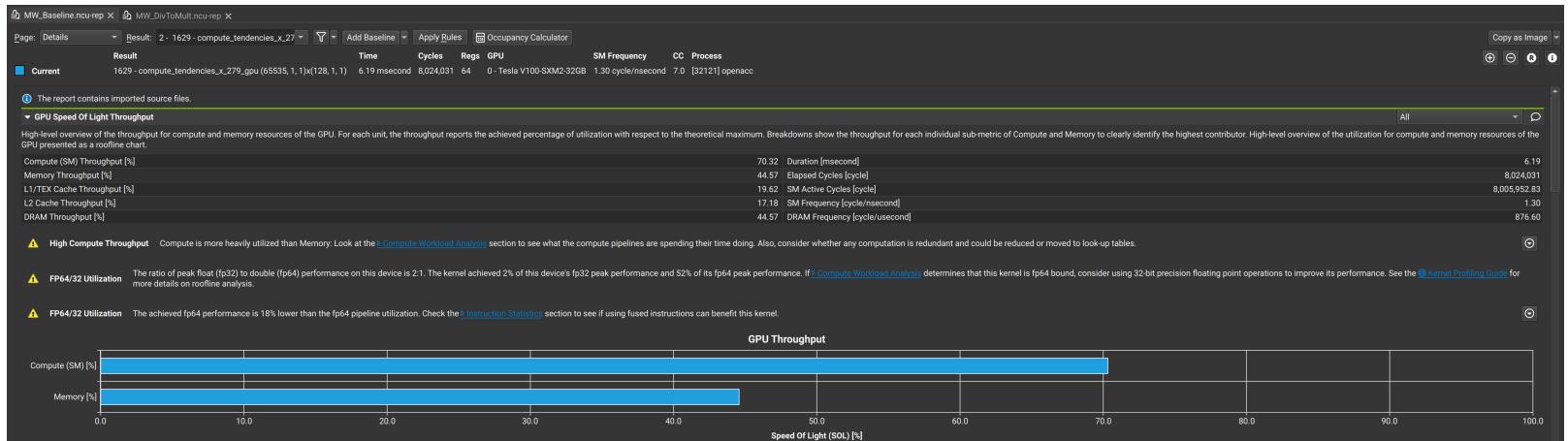
Nsight Compute - Source/SASS and Instruction Hotspots



Navigated to via the **Source Counters** section, a heatmap of resource usage and other metrics can be correlated to specific lines of code within the source files. This can more easily identify which specific areas of your program are causing poor performance.

Nsight Compute - Add a Baseline

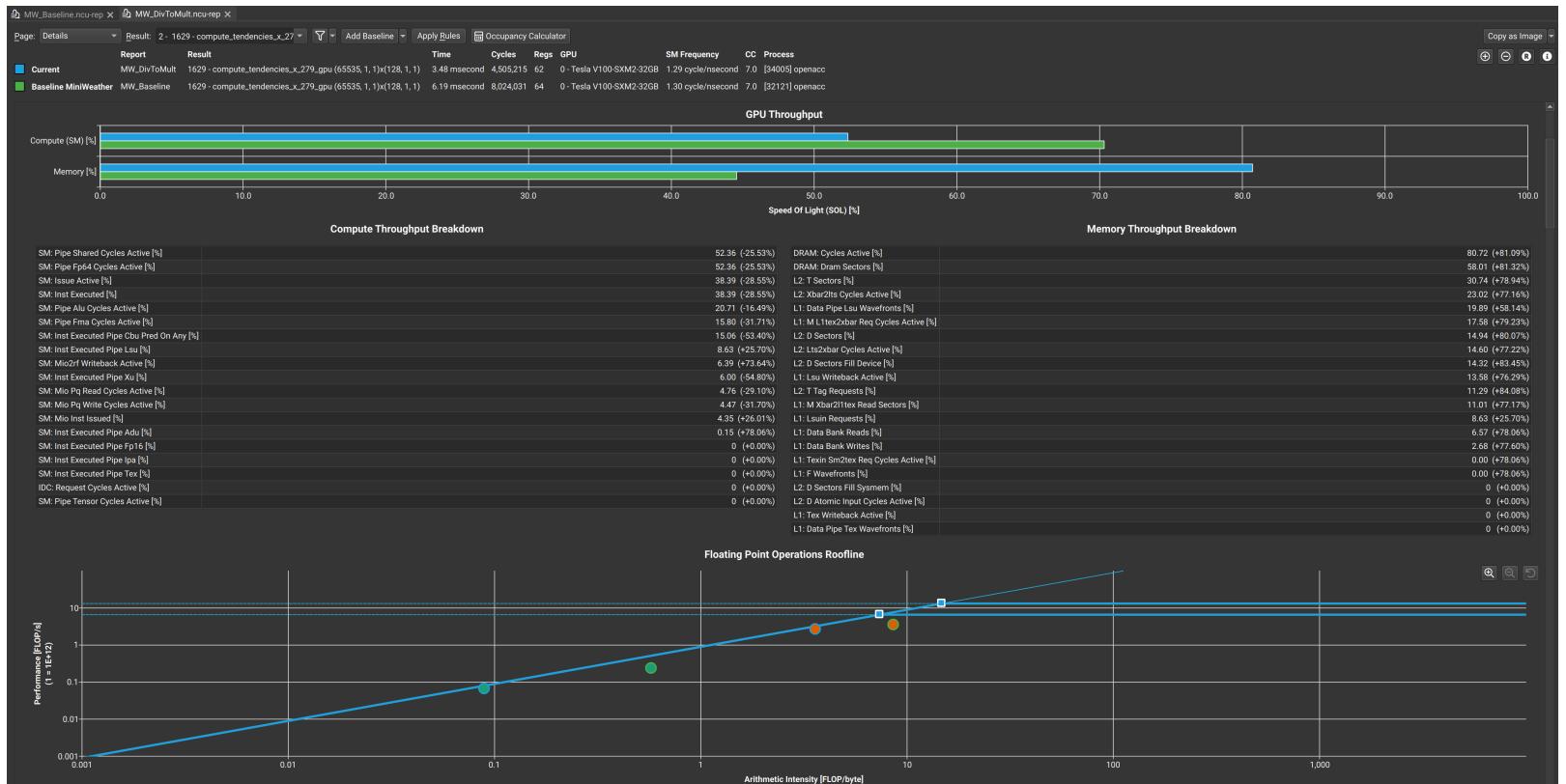
Whenever profiling a program or specific kernel, it is vitally important to record and **set a baseline** to reference performance changes against. In Nsight Compute, set a baseline by clicking **Add Baseline** near the top of the main window within the Nsight Compute GUI. Note, you can add multiple "baselines" from multiple reports.



Rename a baseline by clicking the **Baseline #** text label.



Now, open the new profile report or switch to the other tab referencing this report. The baseline performance metrics will now be displayed and compared to the new current report's performance metrics.



Experiment with a Proposed Optimization - Replace Divide with Multiply

Noting the **hotspot at line 288**, we can assess if there's a way to re-formulate this line to either reduce redundant operations or refactor the overall algorithm. The metrics provided may be able to provide a hint towards why this line is a bottleneck for MiniWeather.

In this case, there are a significant number of warp stalls (see [here](https://docs.nvidia.com/nsight-compute/ProfilingGuide/index.html#statistical-sampler) (<https://docs.nvidia.com/nsight-compute/ProfilingGuide/index.html#statistical-sampler>), for descriptions of types of warp stalls) as well as a much higher number of instructions executed compared to other lines in this kernel. Looking at this line, we see multiple divisions by 12 that could be simplified. Additionally, division is typically more expensive than multiplication within IEEE computational arithmetic.

Thus, let's try changing this line to `vals(ll) = (-stencil(1) + 7*stencil(2) + 7*stencil(3) - stencil(4))*0.0833333333333333`. The report that analyses this change is [MW_DivToMult.ncu-report](#) ([MW_DivToMult.ncu-report](#)).

MW_Baseline.ncu.rep x MW_DivToMult.ncu.rep x

Page: Source Result 2 - 1629 - compute_tendencies_x_27 - Add Baseline Apply Rules Occupancy Calculator

Report Result Time Cycles Regs GPU SM Frequency CC Process

Current MW_DivToMult 1629 - compute_tendencies_x_27.gpu (65535,1,1)(128,1,1) 3.48 msecound 4,505,215 62 0-Tesla V100-SXM2-32GB 1.29 cycle/nsecond 7.0 [34005] openacc

Baseline MiniWeather MW_Baseline 1629 - compute_tendencies_x_27.gpu (65535,1,1)(128,1,1) 6.19 msecound 8,024,031 64 0-Tesla V100-SXM2-32GB 1.30 cycle/nsecond 7.0 [32121] openacc

View: Source and SASI

Source: minnWeather_mpi_openacc.F90 Find... Navigation: L2 Theoretical Sectors Global Excessive

Source Live and Stall Sampling and Stall Sampling Instructions

Registers (All Cycles) (Issued Cycles) Executed

277 hv_coeff = -hv_beta * dx / (16*dt)

278 !Compute fluxes in the x-direction for each cell

279 !\$acc parallel loop collapse(2) private(stencil,vals,d3_vals) async

280 do k = 1 , nz

281 do i = 1 , nx+1

282 !Use fourth-order interpolation from four cell averages to compute the value at

283 do ll = 1 , NUM_VARS

284 do s = 1 , stem_size

285 stencil(s) = state(i-hs-1+s,k,ll)

286 enddo

287 !Fourth-order-accurate interpolation of the state

288 vals(ll) = +stencil(1) + 7*stencil(2) + 7*stencil(3) - stencil(4))*0.088

289 !First-order-accurate interpolation of the third spatial derivative of the state

290 d3_vals(ll) = -stencil(1) + 3*stencil(2) - 3*stencil(1) + stencil(-1)

291 enddo

292

293 !Compute density, u-wind, w-wind, potential temperature, and pressure (r,u,w,r)

294 r = vals(ID_DENs) + hy_dens_cell(k)

295 u = vals(ID_UWOM) / r

296 w = vals(ID_WWOM) / r

297 t = (vals(ID_RHOT) + hy_dens_theta_cell(k)) / r

298 p = C0*(r*t)**gamma

299

300 !Compute the flux vector

301 flux(i,k,ID_DENs) = r*xu - hy_coeff*d3_vals(ID_DENs)

302 flux(i,k,ID_UWOM) = r*xu*p - hy_coeff*d3_vals(ID_UWOM)

303 flux(i,k,ID_WWOM) = r*xw - hy_coeff*d3_vals(ID_WWOM)

304 flux(i,k,ID_RHOT) = r*xt - hy_coeff*d3_vals(ID_RHOT)

305 enddo

306 enddo

307

308 !Use the fluxes to compute tendencies for each cell

309 !\$acc parallel loop collapse(2) async

310 do ll = 1 , NUM_VARS

311 do k = 1 , nz

312 do i = 1 , nx

313 tend(i,k,ll) = -(flux(i+1,k,ll) - flux(i,k,ll)) / dx

314 enddo

315 enddo

316 enddo

317 end subroutine compute_tendencies_x

318

319

320 !Compute the time tendencies of the fluid state using forcing in the z-direction

321 !Since the halos are set in a separate routine, this will not require MPI

322 !First, compute the flux vector at each cell interface in the z-direction (including

Source: compute_tendencies_x_27.gpu Find... Navigation: L2 Theoretical Sectors Global Excessive

Address Source Instructions

126 0008201f fcfae900 LD_E, E, 64, CONSTANT, SYS R28, [R28]

127 0008201f fcfae900 IMAD, X, R25, R27, 0x1, R29, P1

128 0008201f fcfae9f0 LD_E, E, 64, CONSTANT, SYS R10, [R36]

129 0008201f fcfae900 IADD, R28, P1, R30, R1, R2

130 0008201f fcfae910 IMAD, X, R31, R39, 0x1, R29, R9

131 0008201f fcfae920 IADD, R3, P1, R3, c1[0x0][0x100], R2

132 0008201f fcfae930 LD_E, E, 64, CONSTANT, SYS R14, [R14]

133 0008201f fcfae940 IMAD, X, R29, R31, 0x1, R29, P1

134 0008201f fcfae950 LD_E, E, 64, CONSTANT, SYS R24, [R24]

135 0008201f fcfae960 IADD, X, R0, R0, c1[0x0][0x100], R2, P0, IPT

136 0008201f fcfae970 LEA R34, P0, R0, c1[0x0][0x100], R3, R5

137 0008201f fcfae980 LD_E, E, 64, CONSTANT, SYS R18, [R18]

138 0008201f fcfae990 LEA, H, X, R35, R3, c1[0x0][0x174], R8, 0x3, P0

139 0008201f fcfaea00 LD_E, E, 64, CONSTANT, SYS R28, [R28]

140 0008201f fcfaea00 LD_E, E, 64, CONSTANT, SYS R12, [R12]

141 0008201f fcfaea00 LD_E, E, 64, CONSTANT, SYS R34, [R34++0x8]

142 0008201f fcfaead0 LD_E, E, 64, CONSTANT, SYS R8, [R8]

143 0008201f fcfaead0 LD_E, E, 64, CONSTANT, SYS R26, [R26]

144 0008201f fcfaead0 LD_E, E, 64, CONSTANT, SYS R38, [R38]

145 0008201f fcfaeb00 LD_E, E, 64, CONSTANT, SYS R30, [R30]

146 0008201f fcfaeb10 BMOV, Z, CLEAR, R2, B6

147 0008201f fcfaeb20 BSBY, B6, 0xb0fffffae30

148 0008201f fcfaeb30 DFMA R36, R4, _R4

149 0008201f fcfaeb40 DFMA R36, R32, _R4

150 0008201f fcfaeb50 DADD R36, R34, _R54

151 0008201f fcfaeb60 DFMA R22, R36, c1[0x0][0x0], R22

152 0008201f fcfaeb70 MUFO, PCP4H, R37, R23

153 0008201f fcfaeb80 IMAD, MOV, U32, R36, R2, R0, R1

154 0008201f fcfaeb90 DFMA R40, R20, _R4

155 0008201f fcfaebab0 DFMA R42, -R22, R36, 1

156 0008201f fcfaebbb0 DFMA R40, R24, _R4

157 0008201f fcfaebc0 DFMA R42, R42, R42, R42

158 0008201f fcfaebd0 DADD R40, R48, _R28

159 0008201f fcfaebd0 DFMA R44, R44, _R4

160 0008201f fcfaebf0 DFMA R42, R36, R62, R56

161 0008201f fcfaec00 DFMA R4, R46, c1[0x0][0x0], R34

162 0008201f fcfaec10 DFMA R44, R32, -3, R44

163 0008201f fcfaec20 DMUL R58, R4, R42

164 0008201f fcfaec30 DFMA R32, R7, _R14

165 0008201f fcfaec40 DFMA R14, R6, 5, R14

166 0008201f fcfaec50 DFMA R6, R12, 7, _R18

167 0008201f fcfaec60 DFMA R12, R12, 3, -R18

168 0008201f fcfaec70 DFMA R18, -R22, R58, R4

169 0008201f fcfaec80 DFMA R58, R42, R18, R58

170 0008201f fcfaec90 DFMA R10, R28, 3, -R18

171 0008201f fcfaec90 DFMA R32, R6, 7, R32

Copy as Image

EXERCISE - Adjust MiniWeather Problem Size and Other Optimizations

Adjust MiniWeather's problem size using the values

`nx=128, 512, 1024, 2048, 4096, 9192` with `nz=nx/2`. Try more problem sizes if interested. Generate `ncu` reports for each of these problem sizes.

Then, open up all the reports and add each one as a named baseline for that problem size. Compare performance between problem sizes.

- 1. Describe the performance for small problem sizes? What is the SM utilization and memory throughput for small problems?**
- 2. Is there an optimal problem size?**
- 3. Do performance or other metrics stop changing after a certain order of magnitude for the problem size?**
- 4. Experiment with and attempt other optimizations/code changes to improve MiniWeather's performance. What other ways or styles of refactoring might you try to improve performance?**

Resources

- [Nsight Systems Main Documentation](https://docs.nvidia.com/nsight-systems) (<https://docs.nvidia.com/nsight-systems>)
- [Nsight Compute Main Documentation](https://docs.nvidia.com/nsight-compute/) (<https://docs.nvidia.com/nsight-compute/>)
- [Nsight Compute Profiling Guide](https://docs.nvidia.com/nsight-compute/ProfilingGuide/index.html) (<https://docs.nvidia.com/nsight-compute/ProfilingGuide/index.html>)
- [Nsight Compute Training Resources](https://docs.nvidia.com/nsight-compute/Training/index.html) (<https://docs.nvidia.com/nsight-compute/Training/index.html>) - Forum, Videos, and Blog Posts curated by NVIDIA
- Introduction to Kernel Performance Analysis with NVIDIA Nsight Compute, Max Katz (NVIDIA invited to Argonne/NERSC)
 - [GitLab repo](https://gitlab.com/NERSC/roofline-on-nvidia-gpus) (<https://gitlab.com/NERSC/roofline-on-nvidia-gpus>) and [video](https://www.youtube.com/watch?v=fsC3QeZHM1U) (<https://www.youtube.com/watch?v=fsC3QeZHM1U>)

Return to Nsight Systems Profiler Tool

[Nsight Systems Profiler \(./nsys/10_HandsOnNsight_nsys.ipynb\)](#).