

Porting the WAVEWATCH III (v6.07) Wave Action Source Terms to GPU

Olawale James Ikuyajolu^{1,2}, Luke Van Roekel³, Steven R Brus⁴, Erin E Thomas³, Yi Deng^{1,2}, and Sarat Sreepathi⁵

¹Earth and Atmospheric Sciences, Georgia Institute of Technology, Atlanta, GA, USA

²Program in Ocean Science and Engineering, Georgia Institute of Technology, Atlanta, Georgia, USA

³Fluid Dynamics and Solid Mechanics (T-3), Los Alamos National Laboratory, Los Alamos, NM, USA

⁴Mathematics and Computer Science Division, Argonne National Laboratory, Lemont, IL, USA

⁵Computational Sciences and Engineering Division, Oak Ridge National Laboratory, Oak Ridge, TN, USA

Correspondence: Olawale James Ikuyajolu (oikuyajolu3@gatech.edu)

Abstract. Surface gravity waves play a critical role in several processes, including mixing, coastal inundation and surface fluxes. Despite the growing literature on the importance of ocean surface waves, wind-wave processes have traditionally been excluded from Earth system models due to the high computational costs of running spectral wave models. The Next Generation Ocean Model Development in the DOE's (Department of Energy) E3SM (Energy Exascale Earth System Model) project partly focuses on the inclusion of a wave model, WAVEWATCH III (WW3), into E3SM. WW3, which was originally developed for operational wave forecasting, needs to be computationally less expensive before it can be integrated into ESMs. To accomplish this, we take advantage of heterogeneous architectures at DOE leadership computing facilities and the increasing computing power of general-purpose graphics processing units (GPU). This paper identifies the wave action source terms, `w3srcemd`, as the most computationally intensive module in WW3 and then accelerates them via GPU. Our experiments on two computing platforms, Kodiak (P100 GPU & Intel(R) Xeon(R) CPU E5-2695 v4) and Summit (V100 GPU & IBM POWER9) show average speedups of 2x and 4x respectively when mapping one MPI per GPU. An average speedup of 1.4x was achieved using all the 42 CPU cores and 6 GPUs on a summit node (with 7 MPI ranks per GPU). However, the GPU speedup over the 42 CPU cores remains relatively unchanged ($\sim 1.3x$) even when using 4 MPI ranks per GPU (24 ranks in total) and 3 MPI ranks per GPU (18 ranks in total). This corresponds to a 35-40% decrease in both simulation time and usage of resources. Due to too many local scalars and arrays in the `w3srcemd` subroutine and the huge WW3 memory requirement, GPU performance is currently limited by the data transfer bandwidth between the CPU and the GPU. Ideally, OpenACC routine directives could be used to further improve performance. However, `w3srcemd` would require significant code refactoring to make this possible. We also discuss how the trade off between the occupancy, register and latency affects the GPU performance of WW3.

1 Introduction

Ocean surface gravity waves, which derive energy and momentum from steady winds blowing over the surface of the ocean, are a very crucial aspect of the physical processes at the atmosphere-ocean interface. They influence a variety of physical processes such as momentum and energy fluxes, gas fluxes, upper ocean mixing, sea spray production, ice fracture in the marginal ice zone and Earth albedo (Cavaleri et al., 2012). Such complex wave processes can only be treated accurately by including a wave model into Earth system models (ESMs). The first ocean models neglected the existence of ocean waves by assuming that the ocean surface is rigid to momentum and buoyancy fluxes from the atmospheric boundary layer (Bryan and Cox, 1967). Currently, most state-of-the-science ESMs are still missing some or all of these wave-induced effects (Qiao et al., 2013) despite the growing literature on their importance in the simulation of weather and climate.

Recent literature has shown that incorporating different aspects of surface waves into ESMs lead to improved skill performance, particularly in the simulation of sea surface temperature, wind speed at 10m height, ocean heat content, mixed layer depth, and the Walker and Hadley circulations (Law Chune and Aouf, 2018; Song et al., 2012; Shimura et al., 2017; Qiao et al., 2013; Fan and Griffies, 2014; Li et al., 2016). Yet, only two climate models that participated in the Climate Model Intercomparison Project phase 6 (CMIP6), i.e. the First Institute of Oceanography-Earth System Model version 2 (FIO-ESM v2.0; Bao et al., 2020) and the Community Earth System Model Version 2 (CESM2; Danabasoglu et al., 2020) have a wave model as part of their default model components. However, for CMIP6, only FIO-ESM v2.0 employed a wave model. Wind-wave induced physical processes have traditionally been excluded from ESMs due to the high computational cost of running spectral wave models on global model grids for long term climate integrations. In addition to higher computing costs due to longer simulation times, adding new model components also increases resource requirements e.g. number of CPUs/nodes. The Next Generation Ocean Model Development in the US DOE’s (Department of Energy) E3SM (Energy Exascale Earth System Model) project, partly focuses on the inclusion of a spectral wave model, WAVEWATCH III (WW3), into the E3SM to improve the simulation of coastal processes within E3SM. To make WW3 within E3SM feasible for long term global integrations, we need to make it computationally less expensive.

Computer architectures are evolving rapidly, especially in the high-performance computing environment, from traditional homogeneous machines with multicore central processing units (CPUs) to heterogeneous machines with multi-node accelerators such as graphics processing units (GPUs) and multicore CPUs. Moreover, the number of CPU-GPU heterogeneous machines in top 10 of the TOP500 supercomputer increases from 2 in November 2015 (<https://www.top500.org/lists/top500/2015/11/>) to 7 in November 2022 (<https://www.top500.org/lists/top500/2022/11/>). The advent of heterogeneous super-computing platforms, together with the increasing computing power and low energy to performance ratio of GPUs, has motivated the use of GPUs to accelerate climate and weather models. In recent years, numerous studies have reported successful GPU porting of full or partial weather and climate models with improved performances (Hanappe et al., 2011; Xu et al., 2015; Yuan et al., 2020; Zhang et al.,

2020; Bieringer et al., 2021; Mielikainen et al., 2011; Michalakes and Vachharajani, 2008; Shimokawabe et al., 2010; Govett et al., 2017; Li and Van Roekel, 2021; Xiao et al., 2013; Norman et al., 2017, Norman et al., 2022; Bertagna et al., 2020). GPU programming model is different than CPU code, so programmers must recode or use directives to port codes to GPU. Because climate and weather models consist of million lines of code, a majority of the
60 GPU-based climate simulations only operate on certain hot spots (most computationally operations) of the model while leaving a large portion of the model on CPUs. In recent years, however, efforts have been made to run an entire model component on the GPU. For example, Xu et al. (2015) port the entire Princeton Ocean Model to GPU and achieved a 6.9x speedup. Similarly, the entire E3SM atmosphere including SCREAM (Simple Cloud-Resolving E3SM Atmosphere Model) is running on the GPU (<https://github.com/E3SM-Project/scream>). Taking advantage
65 of the recent advancements of GPU programming in climate sciences together with the heterogeneous architectures at DOE leadership computing facilities, this study seeks to identify and move the computationally intensive parts of WW3 to GPU through the use of OpenACC pragmas.

The rest of the paper is structured as follows. In Section 2, we first present an overview of the WW3 model and its parallelization techniques, give an introduction to the OpenACC programming model, describe the hardware
70 and software environment of our testing platforms, and finally present the test case configuration used in this study. The results section, Section 3, presents WW3 profiling analysis on CPU, discusses the challenges encountered, GPU-specific optimization techniques and compares the GPU results with the original FORTRAN code. Section 4 concludes the paper.

2 Model Description and Porting Methodology

75 2.1 WAVEWATCH III

WW3 is a third-generation spectral wave model developed at the US National Centers for Environmental Prediction (NOAA/NCEP) (WAVEWATCH III® Development Group, 2019) from the WAVE Model (WAM) (The Wamdi Group, 1988). It has been used widely to simulate ocean waves in many oceanic regions for various science and engineering applications (Chawla et al., 2013b; Alves et al., 2014; Cornett, 2008; Wang and Oey, 2008). To propagate
80 waves, WW3 solves the random phase spectral action density balance equation, $N(\phi, \lambda, \sigma, \theta, t)$, for wavenumber-direction spectra. The intrinsic frequency (σ) relates the action density spectrum to the energy spectrum (F), $N = \frac{F}{\sigma}$. For large-scale applications, the evolution of the wave action density in WW3 is expressed in spherical coordinates as follows:

$$\frac{\partial N}{\partial t} + \frac{\partial(C_\phi N)}{\partial \phi} + \frac{\partial(C_\lambda N)}{\partial \lambda} + \frac{\partial(C_\sigma N)}{\partial \sigma} + \frac{\partial(C_\theta N)}{\partial \theta} = \sum_i S_i \quad (1)$$

85 Eqn. 1 is solved by discretizing in both physical space (λ, ϕ) and spectral space (σ, θ). Where F is the Energy density; ϕ is the longitude; λ is the latitude; σ is the relative frequency; θ is the direction; t is the time and S

represents the source and sinks terms. The net source-sink terms consist of several physical processes responsible for generation, dissipation and redistribution of energy. The net source-sink terms available in WW3 are waves generation due to wind (S_{in}), dissipation (S_{ds}), non-linear quadruplet interactions (S_{nl}), bottom friction (S_{bt}), and depth-limited breaking (S_{db}), Triad wave-wave interactions (S_{tr}), scattering of waves by bottom features (S_{sc}), wave-ice interactions (S_{ice}) and reflection off shorelines or floating objects (S_{ref}). Details of each source term can be found in the WW3 manual (WAVEWATCH III® Development Group, 2019). The primary source-sink terms used in this work are S_{in} , S_{ds} , S_{nl} , S_{bt} and S_{db} . Several modules are used for the calculation of source terms. However, module `w3srcemd` manages the general calculation and integration source terms.

When moving code between different architectures, it is necessary to understand its structure. For the purposes of our study, Fig. 1 shows a representation of the WW3 algorithm structure. WW3 is divided into several submodules, but the actual wave model is the `w3wavemd`, which runs the wave model for a given time interval. Within `w3wavemd`, several modules are called at each time interval to handle initializations, interpolation of winds and currents, spatial propagation, intra-spectral propagation, calculation and integration of source terms, output file processing, etc. In our work we found that `w3srcemd` is the most computationally intensive part of WW3, thus we focus our attention on this module. According to Fig. 1a, `w3srcemd` is being called at each spatial grid point, which implies that the spatial grid loop is not contained in `w3srcemd`, but rather in `w3wavemd`. Fig. 1b, which represents `w3srcemd`, contains a dynamic integration time loop that can only be executed sequentially. It also calls a number of submodules, such as `w3src4md` for computation of the wind input and wave breaking dissipation source terms. `w3srcemd` consists of collapsed spectral loops ($NSPECH = NK \times NTH$), where NK is the number of frequencies (σ) and NTH is the number of wave directions (θ). Lastly, the `w3src4md` Fig. 1c consists of only frequency (NK) loops. The structure of other source terms submodules is similar to `w3src4md`.

2.2 WW3 Grids and Parallel Concepts

The current version of WW3 can be run and compiled for both single and multi-processor (MPI) compute environments with regular grid, two-way nested (mosaic) grids (Tolman, 2008), spherical multi-cell (SMC) grids (Li, 2012), and unstructured triangular meshes (Roland, 2008, Brus et al., 2021). In this study, we ran and compiled WW3 using MPI with unstructured triangular meshes as the grid configuration. In WW3, the unstructured grid can be parallelized in physical space using either Card Deck (CD) (Tolman, 2002) or domain decomposition (Abdolali et al., 2020). Following Brus et al. (2021), we used the CD approach as the parallelization strategy. The ocean (active) grid cells are sorted and distributed linearly between processors in a round-robin fashion using $n = \text{mod}(m - 1, N)$ i.e. grid cell m is assigned to processor n . Where N is the total number of processors and M is the total number of ocean grids. If N is divisible by M , every processor n has the same number of grids, NSEAL (Fig. 1a). The source term calculation as well as the intra-spectral propagation are computed using the aforementioned parallel strategy, but data are gathered on a single processor to perform the spatial propagation.

To demonstrate the promise of GPU computing for WW3, we used the OpenACC programming model. OpenACC is a directive-based parallel programming model developed to run codes on accelerators without significant programming effort. Programmers incorporate compiler directives in the form of comments into FORTRAN, C, or C++ source codes to assign the computationally intensive the sections of the code to be executed on the accelerator. OpenACC helps to simplify GPU programming because the programmer is not preoccupied with the code parallelism details, unlike CUDA and OpenCL where you need to change the code structure to achieve GPU compatibility. OpenACC compiler automatically transfer calculations and data between two different architectures, the host (CPU) and the accelerator device (GPU). OpenACC works together with OpenMP, MPI and CUDA, supporting heterogeneous parallel environments. Starting from version 4.5, the OpenMP API (Application Programming Interface) specification has been extended to include GPU offloading and GPU parallel directives. While OpenMP and OpenACC have similar constructs, OpenMP is more prescriptive. Prescriptive directives describe the exact computation that should be performed and provide the compiler no flexibility. Our study focuses solely on OpenACC since it has the most mature implementation using the NVIDIA compiler on NVIDIA GPUs at the time of analysis.

OpenACC has three-levels of parallelism (Fig. 2) namely: vector, worker and gangs, corresponding to `threadIdx.x`, `threadIdx.y` and `blockIdx.x` in CUDA terminology . A gang is a group of workers, where multiple gangs work independently without synchronization. Workers are groups of vector/threads within a gang and vector is the finest level of parallelism operating with single-instruction, multiple thread (SIMT). Gang, worker, and vector can be added to a loop region needed to be executed on GPU. An example of Fortran code with and without OpenACC directives is shown in Fig. 3. The OpenACC directives are shown in green as comments starting with `!$acc` (e.g. lines 6 & 12). In Fig. 3a, line 6 is a declaration directive for allocating memory for variables on GPU, line 8 is the data region to move data into the GPU, line 11 updates data already present on GPU with new values from the CPU, line 12 launches the parallel region at the gang level and then dispatches the parallel threads at the worker and vector levels, line 8 updates CPU data with new values from the GPU and line 28 deletes the data on GPU after computation. To learn more about all OpenACC directives, please refer to the NVIDIA (2017) or Chandrasekaran and Juckeland (2017)

2.4 Test Case Configuration

In this study, the WW3 model was configured and simulated over the global ocean with an unstructured mesh of 1°global resolution and 0.25°in regions with depth less than 4km e.g. 1°at the equator and 0.25°at the coastal regions (Brus et al., 2021). The number of the unstructured mesh nodes is 59,072 (ranges from 1°- equator to 0.5°- coastlines; hereafter 59K). In addition, we demonstrate the effect of scaling the problem size on speedup by using an unstructured mesh with 228,540 nodes (ranges from 0.5°- equator to 0.25°- coastlines; hereafter 228K). Following Chawla et al. (2013a), for both spatial grid configurations we use a spectral grid that has 36 directions

and 50 frequency bands that range exponentially from 0.04 to 0.5 Hz , separated by a factor of 1.1. In WW3, the combinations and types of the source terms in Eqn. 1 depends on the research question being answered. WW3 has several source term packages which can be implemented by activating different switches. However, since the goal of this study is purely computational, we selected the commonly used source-sink terms switches, ST4, DB1, BT1 and NL1. The ST4 switch (Ardhuin et al., 2010) consists of the wind input (S_{in}) and wave breaking dissipation (S_{ds}) source terms, DB1 switch is for the depth-induced breaking (S_{db}) source term, BT1 switch consist of bottom friction (S_{bt}) source term parameterizations and the non-linear quadruplet wave interactions (S_{nl}) are computed in the model using the NL1 switch.

We develop and test our accelerated WW3 code on GPU on two computational platforms with heterogeneous architectures, namely:

1. Kodiak cluster from the Parallel Reconfigurable Observational Environment (PROBE) (Gibson et al., 2013) of Los Alamos National Laboratory. Kodiak has 133 compute nodes. Each node contains Intel(R) Xeon(R) CPU E5-2695 v4 @ 2.10GHz 36 CPU cores and 4 NVIDIA Tesla P100 SXM2 GPGPUs, each with 16 GB of memory.
2. Summit, a high performance computing system at the Oak Ridge National Laboratory at Oak Ridge National Laboratory (ORNL). Summit has 4,608 nodes, each contains 2 IBM POWER9 CPUs and 6 NVIDIA Tesla V100 GPUs, each with 80 streaming multiprocessors. All are connected together with NVIDIA’s high-speed NVLink. Summit is the fastest supercomputer in the US and the second fastest in the world in 2021.

Table 2 shows the configuration of each compute node for both platforms. For a more accurate comparison of CPU and GPU codes, we used the same compiler. On Kodiak, the CPU FORTRAN code was compiled using the flags `-g -O3 -acc`. Similarly, the OpenACC code was compiled with flags `-g -O3 -acc -Minfo=accel -ta=tesla,ptxinfo, maxregcount:n`. Likewise on Summit, the flags for the CPU code are `-g -O2`, and that of the OpenACC are `-g -O2 -acc -ta=tesla,ptxinfo,maxregcount:n -Minfo=accel`. Options `-ta=tesla:ptxinfo,maxregcount:n` are the optimization flags¹ used in this study (Section 3.2). Adding the option `-ta=tesla:ptxinfo` to the compile flags provides information about the amount of shared memory used per kernel (a function that is called by the CPU for execution on the GPU) as well the registers per thread. The flag `-ta=tesla:maxregcount:n`, where “n” is the number of registers, sets the maximum number of registers to use per thread.

As a test case, we performed a 5-hour simulation from 2005/06/01 00:00:00 - 05:00:00 by forcing WW3 with atmospheric winds derived from the US National Center for Atmospheric Research reanalysis (NCAR). We verify the GPU model for correctness using significant wave heights (SWH) from the CPU only simulation.

¹WW3 hangs when the -O3 optimization flag is used

3 Results and Evaluations

In this section, we first describe the performance of WW3 on CPU and its computationally intensive sections. Furthermore, we discuss the challenges encountered, how GPU optimization is done and present the performance result of porting WW3 on GPU. Lastly, we discuss the performance limitation using roofline analysis.

3.1 WW3 Profiling Analysis on CPU

With model optimization, an important step is to find the runtime bottlenecks by measuring the performance of various sections in units of time and operations. In order not to waste time and resources improving the performance of rarely used subroutines, we first need to figure out where WW3 spends most of its time. The technical term for this process is called profiling. For this purpose, we profile WW3 by running the Callgrind profiler from the Valgrind tool and then visualize the output using a KCachegrind tool (Weidendorfer, 2008). An application's performance can be 10 to 50 times slower when profiling it with callgrind, however, the proportions of times remain the same. Figure 4 shows the callgraph obtained by profiling WW3 with 300 MPI ranks. The source term subroutine, `w3srcecmd`, can easily be spotted as the consumer of $\sim 82\%$ of the total execution time and resources. Within the `w3srcecmd` subroutine, the dissipation source term S_{ds} uses more than 40% of the total runtime because it contains numerous time consuming spectral loops. In fact, profiling WW3 with another profiler (not shown), Intel Advisor (Intel Corporation, 2021), specifically highlights the spectral time consuming loops. In WW3, each processor serially runs through its sets of allocated spatial grids (described in Section 2.2) with each containing spectral gridpoints, and `w3srcecmd` has a time-dynamic integration (Fig. 1b) procedure which can not be parallelized. Looping through spectral grids and the time-dynamic integration procedure are plausible reasons why `w3srcecmd` is the bottleneck of WW3. We name WW3 model with GPU accelerated `w3srcecmd` as `WW3-W3SRCEMD.gpu` and `WW3.cpu` as the CPU only version.

Fortunately, `w3srcecmd` does not contain neighboring grid dependencies in both spatial and spectral i.e. no parallel data transfers, `w3srcecmd` can therefore be ported to GPUs with less difficulty. We moved the entire WW3 source terms computation to GPU as shown in Fig. 3b.

3.2 Challenges & Optimization

Once the program hotspot is ported to the GPU, the GPU code needs to be optimized in order to improve its performance. Conventionally, the optimization of GPU codes involves loop optimization (fusion and collapse), data transfer management (CPU to GPU and GPU to CPU), memory management and occupancy. Some of these optimization techniques are interrelated e.g. memory management and occupancy. The WW3 model contains very few collapsible loops, so loop fusing and loop collapsing did not effectively optimize the code (figure not shown). To successfully port `WW3-W3SRCEMD.gpu` and achieve the best performance, two challenges had to be overcome in this study. The first is a data transfer issue caused by the WW3 data structure, while the second is a memory

management and occupancy issue caused by the use of many local arrays, sometimes of spectral length, and scalars
215 within `W3SRCEMD` and its embedded subroutines.

3.2.1 Data Transfer Management

It is important to understand the layout of data structures in the program before porting to GPU. WW3 outlines
its data structures by using modules e.g. `w3adatmd`, `w3gdatmd`, & `w3odatmd` (lines 2 & 3 of Fig. 3). Depending on the
variable required, each subroutine uses these modules. These variables are called global external variables. Since
220 it is not possible to move data within a compute kernel in GPU, all necessary data must be present on the GPU
prior to launching the kernel that calls `W3SRCEMD`. The structure of WW3 requires the use of routine directive
(`!$acc routine`) to create a device version of `W3SRCEMD`, as well as other subroutines in it. In FORTRAN, the routine
directive appears in the subprogram specification section or its interface block. Due to the use of routine directives,
OpenACC declare directives(`!$acc declare create`; line 6 of Fig. 3b) were added to the data structures module to
225 inform the compiler that global variables need to be created in the device memory as well as the host memory.

All data must be allocated on the host before being created on a device unless it has been declared device resident.
In WW3, however, arrays whose size are determined by spatial-spectral grid information are allocated at runtime
rather than within the data structure modules. Thus, creating data on the device within the WW3 data modules
poses a problem. Due to this restriction, it was then necessary to explicitly move all the required data to GPU before
230 launching the kernel. For managing data transfers between iteration cycles, we use `!$acc update device(variables)` and
`!$acc update self(variables)` (lines 11 & 21 in Fig. 3b). However, it would be easier to use unified memory, which
offers a unified memory address space to both CPUs and GPUs, rather than tracking each and every data that needs
to be sent to the GPU, but the latest OpenACC version does not support the use of unified memory with routine
directives. In the future, having this feature would save time spent tracking data transfers in programs with many
235 variables such as WW3.

3.2.2 Memory management and Occupancy

Occupancy is defined as the ratio of active warps (workers) on a streaming multiprocessor (SM) to the maximum
number of active warps supported by the SM. On Kodiak and Summit, the maximum number of threads per SM
is 2048. A warp consists of 32 threads which implies that the total number of possible warps per SM is 8. Even
240 if the kernel launch configuration maximizes the number of threads per SM, other GPU resources, such as shared
memory and registers, may also limit the number of maximum threads, thus indirectly affecting GPU occupancy. A
register is a small amount of fast storage available to each thread, whereas shared memory is the memory shared by
all threads in each SM. As seen in Table 2, the maximum memory per SM for Kodiak and Summit is 256KB, but
the maximum shared memory for Kodiak is 64KB (P100) and 94KB for Summit (V100). In terms of total register
245 file size per GPU, Kodiak and Summit have 14336KB and 20480KB, respectively.

To estimate GPU usage, we use the CUDA occupancy calculator available on https://docs.nvidia.com/cuda/cuda-occupancy-calculator/CUDA_Occupancy_Calculator.xls. In this study, the kernel launch configurations consist of NSEAL gangs, where NSEAL refers to the number of grids on each node, and each gang has 32 vector lengths (or threads). With this configuration, the achievable GPU occupancy (regardless of other resources) is 50%. Adding `-ta=tesla:pxtinfo` to the compiler flags gives the information about the size of registers, memory spills (movement of variables out of the register space to the main memory) and shared memory allocated during compilation. Kodiak and Summit both allocate the maximum 255KB register per thread, reducing GPU occupancy to 13%. In addition, a full register leads to spilling of memory into the L1 cache (shared memory). A spill to cache is fine, but a spill to the global memory will severely affect performance because the time required to get data from the global memory is longer than that from a register, latency. Latency is the amount of time required to move data from one point to the other. Therefore, an increase (decrease) in register size causes two different things simultaneously: a decrease (increase) in occupancy and decrease (increase) in latency. There is always a trade off between register, latency and occupancy, and the the trick is to find the spot that maximizes performance. One can set the maximum number of registers per thread via the flag `-ta=tesla:maxregcount:n` where “n” is the number of registers.

Figure 5 illustrates how the trade off between latency and occupancy affects the GPU performance based on the number of registers. Our analysis was based only on the performance of running 228K mesh configuration with 16 MPI tasks. For Summit, as the number of registers increases from 16 to 64, the GPU occupancy remains constant at 50% and GPU performance improved due to the movement of more variables to the fast memory. As indicated by the blue part of the line, this is a latency-dominant region. However, as the number of register increases from 64 to 192, GPU occupancy gradually decreases from 50% to 13%. This degraded performance despite moving more data into the fast memory. Therefore, this is an occupancy-dominant region as indicated by the red part of the line (Fig. 5b). From 192 to the maximum register count, the occupancy remains constant at 13% and GPU performance remains relatively constant. As occupancy is constant, latency is expected to dominate this region, which is represented by the green part of the line in Fig. 5b. However, we observe no memory spill for registers between 192 and 255 and thus latency effect remains unchanged. Therefore, constant latency and occupancy result in constant performances in this region.

Fig. 5b shows that 64 registers produced the best performance (minimum runtime) on Summit. For brevity, the previously described trade-off analysis also applies to Kodiak, and 96 register count produced the best performance (Fig. 5a). On Kodiak, the register count that achieved the best performance is higher than on Summit, probably due to Summit’s larger L1 and L2 caches (Table 2). With larger L1 and L2 caches, more data can be stored, reducing memory spillover to global memory and thus reducing latency. Optimizing the register count increases the GPU performance by approximately 20% on Kodiak and 14% on Summit.

3.3 GPU Accelerated W3SRCEMD

280 This section compares the performance of WW3.cpu and WW3-W3SRCEMD.gpu. Note that the speedups in this section are achieved by parallelizing the local grids loop, which calls the W3SRCEMD function, using the OpenACC parallel directives (Fig. 3b). In `w3srcecmd` and its dependent subroutines, we introduced the OpenACC routine directive (`!$acc routine`) which instructs the compiler to build a device version of the subroutine so that it may be called from a device region by each gang. In addition, at the start of the time integration, we moved the needed constants data to
285 the GPU (line 8 of Fig. 3b). Using the average over the late 2-hour simulation, Fig. 7 compares the output results of CPU and WW3-W3SRCEMD.gpu codes and their relative difference for significant wave height (SWH). According to the validation results, the SWH is nearly identical, and the error is negligible and acceptable. It is possible that the error stems from the difference in mathematical precision between the GPU and CPU.

For simplicity, we start by mapping one MPI rank to one GPU. Comparing the performance of a single GPU with
290 a single CPU core (Table 2) on Kodiak, a 2.3x & 2.4x speedup was achieved for 59K and 228K meshes respectively. Here, the speedup is relative on a single CPU core. Similarly on Summit, we achieved speedups of 4.7x & 6.6x for 59K and 228K meshes respectively. On Summit, the GPU performance of 228K nodes is better because the CPU gets extremely slow. Summit’s speedup is greater than Kodiak’s because the Tesla V100-SXM2-16GB GPU is faster than the Tesla P100-PCIE-16GB GPU (NVIDIA, 2017). Due to the reduction in GPU workload, speedups gradually
295 decline as the number of MPI ranks increases (Table 2).

3.3.1 Fair Comparison Using Multi-Process Service (MPS)

The following sections focus exclusively on Summit’s results. Each summit compute node is equipped with six GPUs and 42 CPU cores. As an initial step, we used only six MPI ranks on each node so that each process could offload its work to one GPU. In order to properly compare CPUs and GPUs on a single node, we must use the whole
300 CPU and GPU resources on a node. The NVIDIA Volta GPUs support multiprocessing services. When running with a mesh size greater than 59K and assigning multiple MPI ranks to a single GPU, the GPU heapsize overflows with an auto allocation problem. The auto allocation problem is solved by setting the environment variable “`PGI_ACC_CUDA_HEAPSIZE`” to a minimum allowed number based on the mesh size and MPI ranks per GPU. By mapping 7 MPI ranks per GPU on summit, a speedup of 1.36 (1.41) was achieved for the 59K (228K) mesh size
305 (Fig. 6) over the whole CPU cores. We also ran 4, 3 and 2 MPI ranks per GPU configurations. The results of the different multi-process configurations and the two mesh sizes are shown in Table 3. Note that the speedups of all the GPU configurations are compared with the full CPU 42 cores run. By varying the number of MPI ranks per GPU, the 4 and 3 MPI ranks per GPU (24 and 18 MPI ranks in total) have higher speedups than 7 MPIs per GPU for the 59K mesh size. Even with 2 MPI ranks per GPU (total of 12 ranks), the speedup is 1.23x. The speedup for
310 the 228K gradually decreases from 1.41x to 1.12x as the number of MPI ranks per GPU decreases from 7 to 2. The workload (distributed grids) per MPI rank reduces as the number of MPI ranks increases. As a result of the reduced

workload, GPU utilization falls. The overall performance of a GPU depends on both the number of MPI ranks per GPU and the number of grid size per MPI rank.

315 Table 4 shows the results of scaling the 228K mesh size over multiple nodes. We used the full GPU configuration here by launching 7 MPI ranks on each GPU. In the results, it can be seen that the speedup is relatively uniform across multiple nodes. Likewise, as the gridpoints per MPI rank decreases with increasing nodes, speedup decreases gradually due to reduced GPU utilization. However, speedup is always the highest whenever the number of gridpoints per MPI rank is \sim around 2500 e.g. speedup of 2 nodes in Table 4 for 228K and the speedup of 24 MPI
320 ranks for 59K mesh in Table 3. This is most likely due to the tail effect - load imbalance between SMs. Based our kernel launch configuration, 64 register and 32 block size, the number of possible blocks per SM is limited to 32. With 80 SM on a V100 GPU, the maximum total number of blocks (gangs i.e. gridpoints) that could be executed simultaneously is 80×32 (2560).

325 All previous speedups are based on the whole WW3 code. However, only `w3srceMD` is being accelerated on the GPU and the rest of the code run on CPU. So when comparing 24 MPI ranks GPU configuration with all the 42 CPU cores, 24 MPI ranks is being used for the CPU section of `WW3-w3srceMD.gpu` against 42 MPI ranks for the `WW3.cpu`. Therefore, it is necessary to also consider only the speedup of `w3srceMD` on GPU over the CPU. The last row of Table 3 (also Fig. 6) shows the runtimes and the achieved speedup of `w3srceMD` subroutine for 288K mesh size
330 on summit node. In comparison with the speedup based on the whole WW3 code, `w3srceMD` speedup increases for all configurations with a maximum speedup of 1.61x.

3.4 W3SRCEMD Roofline Plot

The roofline model helps us understand the trade-off between data-movement and computation, so we can find out
335 what's limiting our code and how close we are to it. The roofline (Fig. 8) indicates that, as expected from a memory-intensive model, the kernel is limited by the data transfer bandwidth between the CPU and the GPU. Most of the kernel's time is spent in executing memory (load/store) instructions. It is worth mentioning that the device achieved a compute throughput and a memory bandwidth that are both below 40% of its peak performance. Therefore, it appears that, while the computation is waiting for the GPU to provide needed data, the GPU is waiting for the
340 CPU to transfer data. Similarly, the roofline shows that the kernel memory bandwidth is approximately equal to the NVLink bandwidth. Thus, the `w3srceMD` kernel is limited by the bandwidth between the CPU and the GPU memory.

Also, the kernel has a very low arithmetic intensity for both simple and double-precision floating-point (Fig. 8) computations, performing only very few flops for every double and integer loaded from DRAM. Arithmetic intensity is a measure of floating-point operations (FLOPs) performed relative to the amount of memory accesses in bytes
345 that are needed for those operations. For our 5-hour simulation, the kernel is being launched 21 times and therefore

42 data movements between the host and the device for non-constant variables. Unfortunately, the large arrays need to be updated on both device and host at each time step. As an example, VA, the spectra storage array, is approximately 5Gb (20Gb) for a spatial mesh size of 59,000 (228,000) and spectral resolution of 50x36. In summary, W3SRCE subroutine is simply too big and complicated to be ported efficiently using OpenACC routine directives
350 and therefore requires significant refactoring.

4 Conclusions

Climate science is increasingly moving towards higher spatial resolution models and additional components to better simulate previously parameterized or excluded processes. In recent decades, the use of GPUs to accelerate scientific problems has increased significantly due to the emergence of supercomputers with heterogeneous architectures.

355 Wind generated waves play an important role in modifying physical processes at the atmosphere ocean interface. They have generally been excluded from most coupled Earth system models partly due to its high computational cost. However, the Energy Exascale Earth System Model (E3SM) project seeks to include a wave model (WW3) and introducing WW3 to E3SM would increase the computational time and usage of resources.

360 In this study, we identified and accelerated the computationally intensive section of WW3 on GPU using OpenACC. Using Valgrind & callgrind tools, we found that the source term subroutine, w3srcemd, consumes 78% of the execution time. The w3srcemd subroutine has no neighboring gridpoints dependencies, and is therefore well suited for implementation on GPU. On two different computational platforms, Kodiak with four P100 GPUs and Summit with six V100 GPUs on each node, we performed 5-hour simulation experiments using two global unstructured meshes
365 with 59,072 and 228,540 nodes. On average, running W3RCE by offloading one MPI per GPU gives approximately 4x (2x) speedup over CPU version on Summit (Kodiak). On a fair comparison, by using all the 42 CPU cores and 6 GPUs on a summit node, a speedup of 1.4x was achieved using 7 MPI ranks per GPU. However, the GPU speedup over the 42 CPU cores remain relatively unchanged (~1.3x) even when using 4 MPI ranks per GPU (24 ranks in total) and 3 MPI ranks per GPU (18 ranks in total). The GPU performance is heavily affected by the
370 data transfer bandwidth between the CPU and the GPU. The large number of local scalars and arrays within the w3srcemd subroutine and the huge size of memory required to run WW3 is currently hurting the GPU utilization, thus the achievable speedup. Too many constants in WW3 occupy the register (fast memory) and then spill over to the L1 and L2 caches or the GPU global memory. To increase the GPU performance, the grid loop counter within W3WAVEMD must be pushed into the w3srcemd, thereby moving the gang level parallelization into w3srcemd. This
375 require major code refactoring starting with modification of WW3 data structures. There are other parts of WW3 code that can be ported to GPU, such as spatial and spectral propagation.

Coupling CESM with WW3 at low resolution, Li et al., 2016 found 36% and 28% increase in computational cost for ocean-wave only and fully coupled simulations when running WW3 on a $3.2^\circ \times 4^\circ$ latitude-longitude grid with 25 frequency and 24 directional bins with 3° resolution ocean model and T31 atmosphere. Likewise, in a one-way coupling of WW3 to E3SM atmospheric component, WW3 increases the number of processors by $\sim 35\%$ and its runtime is 44% more than the ocean model. From our first attempt at GPU-based spectral wave modeling, runtime decreased by 35-40% and resource usage decreased by 40-55%. Thus, leveraging heterogeneous architectures reduces the amount of time and resources required to include WW3 in global climate models. It is important to note that WW3 will become a bottleneck as the ocean and atmosphere models in E3SM move towards heterogeneous architectures. Consequently, WW3 needs huge refactoring to take advantage of GPU capabilities and be fully prepared for Exascale regime.

After refactoring the WW3 code, we also need to investigate how different WW3 setups (grid parallelization method, source term switches, propagation schemes, etc.) affect GPU performance. The performance of GPU-accelerated WW3 using OpenMP should also be considered for future work. The success of this work has laid the foundation for future work in global spectral wave modeling, and it is also a major step toward expanding E3SM's capability to run with waves on heterogeneous architectures in the near future.

Code and data availability. Model configuration and input files can be assessed at <https://doi.org/10.5281/zenodo.6483480>. The official repository of WAVEWATCH III CPU code can be found here: <https://github.com/NOAA-EMC/WW3>. The new WW3-W3SRCEMD.gpu code used in this work is available at <https://doi.org/10.5281/zenodo.6483401>

Author contributions. OJI was responsible for the code modifications, simulations, verification tests and performance analyses. LVR developed the concept for this study. OJI provided the initial draft of the manuscript. LVR, SRB, EET, YD and SS contributed to the final manuscripts.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. This research was supported as part of the Energy Exascale Earth System Model (E3SM) project, funded by the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research. This research used resources provided by the Los Alamos National Laboratory (LANL) Institutional Computing Program, which is supported by the US Department of Energy National Nuclear Security Administration under contract no. 89233218CNA000001. This research used resources from the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is

supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725. We wish to thank Phil Jones at LANL for assistance and comments on early drafts. Lastly, we would like to thank Kim Youngsung at ORNL for his technical assistance on the revised manuscript.

References

- 410 Abdolali, A., Roland, A., van der Westhuysen, A., Meixner, J., Chawla, A., Hesser, T. J., Smith, J. M., and Sikiric, M. D.: Large-scale hurricane modeling using domain decomposition parallelization and implicit scheme implemented in WAVEWATCH III wave model, *Coastal Engineering*, 157, 103 656, <https://doi.org/https://doi.org/10.1016/j.coastaleng.2020.103656>, 2020.
- Alves, J.-H. G. M., Chawla, A., Tolman, H. L., Schwab, D., Lang, G., and Mann, G.: The operational implementation of a
415 Great Lakes wave forecasting system at NOAA/NCEP, *Weather Forecast.*, 29, 1473–1497, 2014.
- Bao, Y., Song, Z., and Qiao, F.: FIO-ESM Version 2.0: Model Description and Evaluation, *Journal of Geophysical Research: Oceans*, 125, e2019JC016 036, <https://doi.org/https://doi.org/10.1029/2019JC016036>, e2019JC016036 2019JC016036, 2020.
- Bertagna, L., Guba, O., Taylor, M. A., Foucar, J. G., Larkin, J., Bradley, A. M., Rajamanickam, S., and Salinger, A. G.:
420 A Performance-Portable Nonhydrostatic Atmospheric Dycore for the Energy Exascale Earth System Model Running at Cloud-Resolving Resolutions, SC '20, IEEE Press, 2020.
- Bieringer, P. E., Piña, A. J., Lorenzetti, D. M., Jonker, H. J. J., Sohn, M. D., Annunzio, A. J., and Fry, R. N.: A Graphics Processing Unit (GPU) Approach to Large Eddy Simulation (LES) for Transport and Contaminant Dispersion, *Atmosphere*, 12, 890, <https://doi.org/10.3390/atmos12070890>, 2021.
- 425 Brus, S. R., Wolfram, P. J., Van Roekel, L. P., and Meixner, J. D.: Unstructured global to coastal wave modeling for the Energy Exascale Earth System Model using WAVEWATCH III version 6.07, *Geoscientific Model Development*, 14, 2917–2938, <https://doi.org/10.5194/gmd-14-2917-2021>, 2021.
- Bryan, K. and Cox, M. D.: A numerical investigation of the oceanic general circulation, *Tellus*, 19, 54–80, <https://doi.org/10.3402/tellusa.v19i1.9761>, 1967.
- 430 Cavaleri, L., Fox-Kemper, B., and Hemer, M.: Wind Waves in the Coupled Climate System, *Bulletin of the American Meteorological Society*, 93, 1651 – 1661, <https://doi.org/10.1175/BAMS-D-11-00170.1>, 2012.
- Chandrasekaran, S. and Juckeland, G.: *OpenACC for Programmers: Concepts and Strategies (1st Ed.)*, Addison-Wesley Professional, 2017.
- Chawla, A., Spindler, D. M., and Tolman, H. L.: Validation of a thirty year wave hindcast using the Climate Forecast System
435 Reanalysis winds, *Ocean Model. (Oxf.)*, 70, 189–206, 2013a.
- Chawla, A., Tolman, H. L., Gerald, V., Spindler, D., Spindler, T., Alves, J.-H. G. M., Cao, D., Hanson, J. L., and Devaliere, E.-M.: A multigrid wave forecasting model: A new paradigm in operational wave forecasting, *Weather Forecast.*, 28, 1057–1078, 2013b.
- Cornett, A. M.: A global wave energy resource assessment, in: *The Eighteenth International Offshore and Polar Engineering Conference*, International Society of Offshore and Polar Engineers, 2008.
- 440 Danabasoglu, G., Lamarque, J.-F., Bacmeister, J., Bailey, D. A., DuVivier, A. K., Edwards, J., Emmons, L. K., Fasullo, J., Garcia, R., Gettelman, A., Hannay, C., Holland, M. M., Large, W. G., Lauritzen, P. H., Lawrence, D. M., Lenaerts, J. T. M., Lindsay, K., Lipscomb, W. H., Mills, M. J., Neale, R., Oleson, K. W., Otto-Bliesner, B., Phillips, A. S., Sacks, W., Tilmes, S., van Kampenhout, L., Vertenstein, M., Bertini, A., Dennis, J., Deser, C., Fischer, C., Fox-Kemper, B.,
445 Kay, J. E., Kinnison, D., Kushner, P. J., Larson, V. E., Long, M. C., Mickelson, S., Moore, J. K., Nienhouse, E., Polvani,

- L., Rasch, P. J., and Strand, W. G.: The Community Earth System Model Version 2 (CESM2), *Journal of Advances in Modeling Earth Systems*, 12, e2019MS001916, <https://doi.org/https://doi.org/10.1029/2019MS001916>, e2019MS001916 2019MS001916, 2020.
- Fan, Y. and Griffies, S. M.: Impacts of Parameterized Langmuir Turbulence and Nonbreaking Wave Mixing in Global Climate Simulations, *Journal of Climate*, 27, 4752 – 4775, <https://doi.org/10.1175/JCLI-D-13-00583.1>, 2014.
- Gibson, G., Grider, G., Jacobson, A., and Lloyd, W.: PRObE: A thousand-node experimental cluster for computer systems research, *Usenix ;login*, 38, 2013.
- Govett, M., Rosinski, J., Middlecoff, J., Henderson, T., Lee, J., MacDonald, A., Wang, N., Madden, P., Schramm, J., and Duarte, A.: Parallelization and Performance of the NIM Weather Model on CPU, GPU, and MIC Processors, *Bulletin of the American Meteorological Society*, 98, 2201 – 2213, <https://doi.org/10.1175/BAMS-D-15-00278.1>, 2017.
- Hanappe, P., Beurivé, A., Laguzet, F., Steels, L., Bellouin, N., Boucher, O., Yamazaki, Y. H., Aina, T., and Allen, M.: FAMOUS, faster: using parallel computing techniques to accelerate the FAMOUS/HadCM3 climate model with a focus on the radiative transfer algorithm, *Geoscientific Model Development*, 4, 835–844, <https://doi.org/10.5194/gmd-4-835-2011>, 2011.
- Intel Corporation : Intel Advisor User Guide Version 2022.0, Intel Corporation, <https://www.intel.com/content/www/us/en/develop/documentation/advisor-user-guide/top.html>, 2021.
- Jiang, J., Lin, P., Wang, J., Liu, H., Chi, X., Hao, H., Wang, Y., Wang, W., and Zhang, L.: Porting LASG/ IAP Climate System Ocean Model to Gpus Using OpenAcc, *IEEE Access*, 7, 154 490–154 501, <https://doi.org/10.1109/ACCESS.2019.2932443>, 2019.
- Law Chune, S. and Aouf, L.: Wave effects in global ocean modeling: parametrizations vs. forcing from a wave model, *Ocean Dynamics*, 68, 1739–1758, <https://doi.org/10.1007/s10236-018-1220-2>, 2018.
- Li, J.-G.: Propagation of ocean surface waves on a spherical multiple-cell grid, *Journal of Computational Physics*, 231, 8262–8277, <https://doi.org/https://doi.org/10.1016/j.jcp.2012.08.007>, 2012.
- Li, Q. and Van Roekel, L.: Towards multiscale modeling of ocean surface turbulent mixing using coupled MPAS-Ocean v6.3 and PALM v5.0, *Geoscientific Model Development*, 14, 2011–2028, <https://doi.org/10.5194/gmd-14-2011-2021>, 2021.
- Li, Q., Webb, A., Fox-Kemper, B., Craig, A., Danabasoglu, G., Large, W. G., and Vertenstein, M.: Langmuir mixing effects on global climate: WAVEWATCH III in CESM, *Ocean Modelling*, 103, 145–160, <https://doi.org/https://doi.org/10.1016/j.ocemod.2015.07.020>, waves and coastal, regional and global processes, 2016.
- Michalakes, J. and Vachharajani, M.: GPU acceleration of numerical weather prediction, in: 2008 IEEE International Symposium on Parallel and Distributed Processing, pp. 1–7, <https://doi.org/10.1109/IPDPS.2008.4536351>, 2008.
- Mielikainen, J., Huang, B., and Huang, H.-L. A.: GPU-Accelerated Multi-Profile Radiative Transfer Model for the Infrared Atmospheric Sounding Interferometer, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 4, 691–700, <https://doi.org/10.1109/JSTARS.2011.2159195>, 2011.
- Norman, M. R., Mametjanov, A., and Taylor, M. A.: Exascale Programming Approaches for the Accelerated Model for Climate and Energy, 2017.
- Norman, M. R., Bader, D. A., Eldred, C., Hannah, W. M., Hillman, B. R., Jones, C. R., Lee, J. M., Leung, L. R., Lyngaas, I., Pressel, K. G., Sreepathi, S., Taylor, M. A., and Yuan, X.: Unprecedented cloud resolution in a GPU-enabled full-physics atmospheric climate simulation on OLCF’s summit supercomputer, *Int. J. High Perform. Comput. Appl.*, 36, 93–105, 2022.

- NVIDIA: NVIDIA Tesla V100 GPU Architecture, Tech. rep., NVIDIA Corporation, available at
485 <http://www.nvidia.com/object/volta-architecture-whitepaper.html>, 2017.
- Qiao, F., Song, Z., Bao, Y., Song, Y., Shu, Q., Huang, C., and Zhao, W.: Development and evaluation of an Earth System Model with surface gravity waves, *Journal of Geophysical Research: Oceans*, 118, 4514–4524, <https://doi.org/https://doi.org/10.1002/jgrc.20327>, 2013.
- Roland, A.: Development of WWM II: Spectral wave modeling on unstructured meshes, Ph.D. thesis, 2008.
- 490 Shimokawabe, T., Aoki, T., Muroi, C., Ishida, J., Kawano, K., Endo, T., Nukada, A., Maruyama, N., and Matsuoka, S.: An 80-Fold Speedup, 15.0 TFlops Full GPU Acceleration of Non-Hydrostatic Weather Model ASUCA Production Code, in: SC '10: Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis, pp. 1–11, <https://doi.org/10.1109/SC.2010.9>, 2010.
- Shimura, T., Mori, N., Takemi, T., and Mizuta, R.: Long-term impacts of ocean wave-dependent roughness on global climate
495 systems, *Journal of Geophysical Research: Oceans*, 122, 1995–2011, <https://doi.org/https://doi.org/10.1002/2016JC012621>, 2017.
- Song, Z., Qiao, F., and Song, Y.: Response of the equatorial basin-wide SST to non-breaking surface wave-induced mixing in a climate model: An amendment to tropical bias, *Journal of Geophysical Research: Oceans*, 117, <https://doi.org/https://doi.org/10.1029/2012JC007931>, 2012.
- 500 The Wamdi Group: The WAM model—A third generation ocean wave prediction model, *J. Phys. Oceanogr.*, 18, 1775–1810, 1988.
- Tolman, H. L.: Distributed-memory concepts in the wave model WAVEWATCH III, *Parallel Computing*, 28, 35–52, [https://doi.org/https://doi.org/10.1016/S0167-8191\(01\)00130-2](https://doi.org/https://doi.org/10.1016/S0167-8191(01)00130-2), 2002.
- Tolman, H. L.: A mosaic approach to wind wave modeling, *Ocean Modelling*, 25, 35–47,
505 <https://doi.org/https://doi.org/10.1016/j.ocemod.2008.06.005>, 2008.
- Wang, D.-P. and Oey, L.-Y.: Hindcast of waves and currents in Hurricane Katrina, *B. Am. Meteorol. Soc.*, 89, 487–496, 2008.
- WAVEWATCH III® Development Group: User manual and system documentation of WAVEWATCH III version 6.07. Tech. Note 333, NOAA/NWS/NCEP/MMAB, Tech. rep., College Park, MD, USA, 2019.
- 510 Weidendorfer, J.: Sequential Performance Analysis with Callgrind and KCachegrind, in: Tools for High Performance Computing, edited by Resch, M., Keller, R., Himmler, V., Krammer, B., and Schulz, A., pp. 93–113, Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.
- Xiao, H., Sun, J., Bian, X., and Dai, Z.: GPU acceleration of the WSM6 cloud microphysics scheme in GRAPES model, *Computers & Geosciences*, 59, 156–162, <https://doi.org/https://doi.org/10.1016/j.cageo.2013.06.016>, 2013.
- 515 Xu, S., Huang, X., Oey, L.-Y., Xu, F., Fu, H., Zhang, Y., and Yang, G.: POM.gpu-v1.0: a GPU-based Princeton Ocean Model, *Geoscientific Model Development*, 8, 2815–2827, <https://doi.org/10.5194/gmd-8-2815-2015>, 2015.
- Yuan, Y., Shi, F., Kirby, J. T., and Yu, F.: FUNWAVE-GPU: Multiple-GPU Acceleration of a Boussinesq-Type Wave Model, *Journal of Advances in Modeling Earth Systems*, 12, e2019MS001957, <https://doi.org/https://doi.org/10.1029/2019MS001957>, e2019MS001957 10.1029/2019MS001957, 2020.
- 520 Zhang, S., Fu, H., Wu, L., Li, Y., Wang, H., Zeng, Y., Duan, X., Wan, W., Wang, L., Zhuang, Y., Meng, H., Xu, K., Xu, P., Gan, L., Liu, Z., Wu, S., Chen, Y., Yu, H., Shi, S., Wang, L., Xu, S., Xue, W., Liu, W., Guo, Q., Zhang, J., Zhu, G.,

Tu, Y., Edwards, J., Baker, A., Yong, J., Yuan, M., Yu, Y., Zhang, Q., Liu, Z., Li, M., Jia, D., Yang, G., Wei, Z., Pan, J., Chang, P., Danabasoglu, G., Yeager, S., Rosenbloom, N., and Guo, Y.: Optimizing high-resolution Community Earth System Model on a heterogeneous many-core supercomputing platform, *Geoscientific Model Development*, 13, 4809–4829, <https://doi.org/10.5194/gmd-13-4809-2020>, 2020.

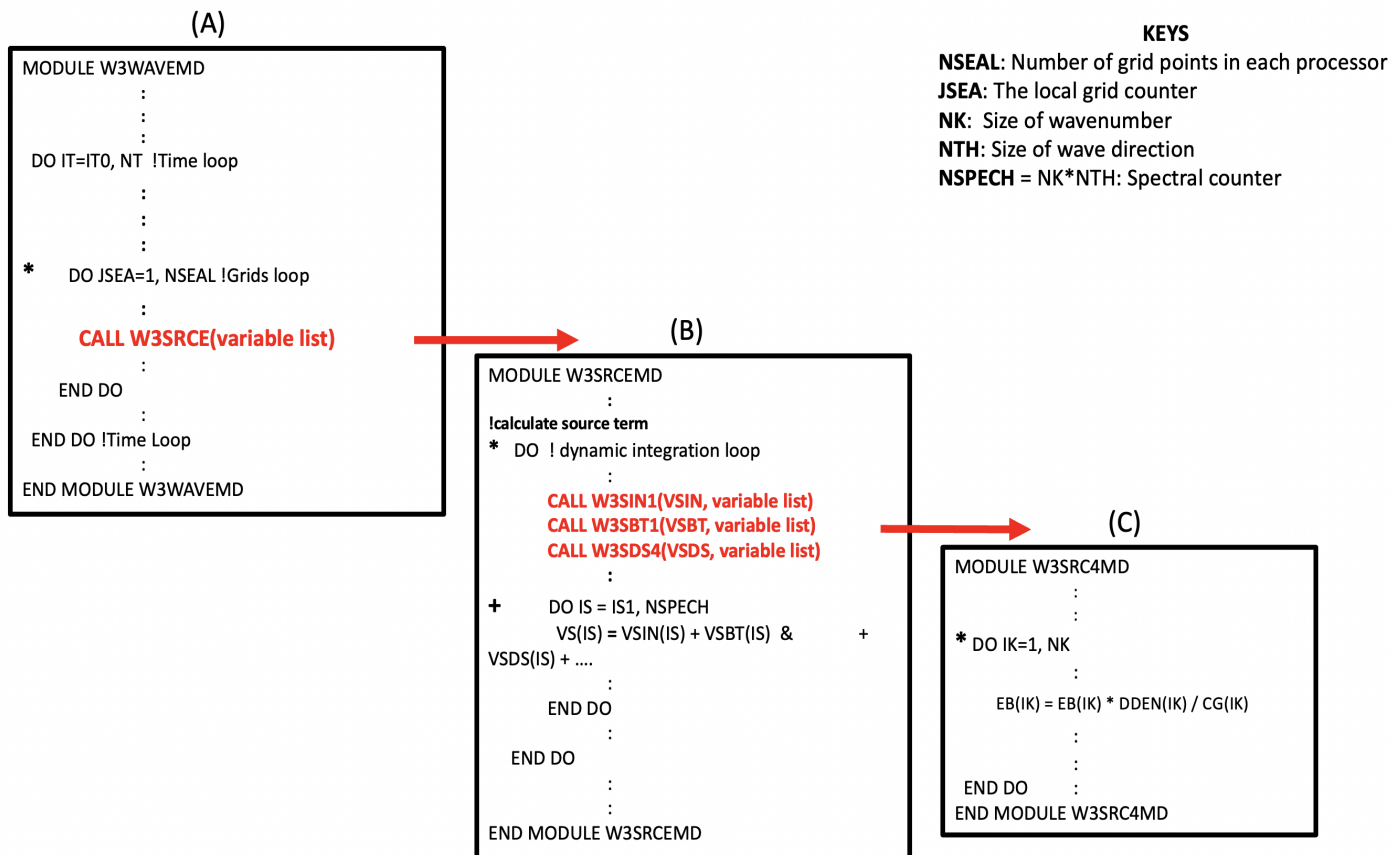


Figure 1. A schematic representation of WAVEWATCH III code structure

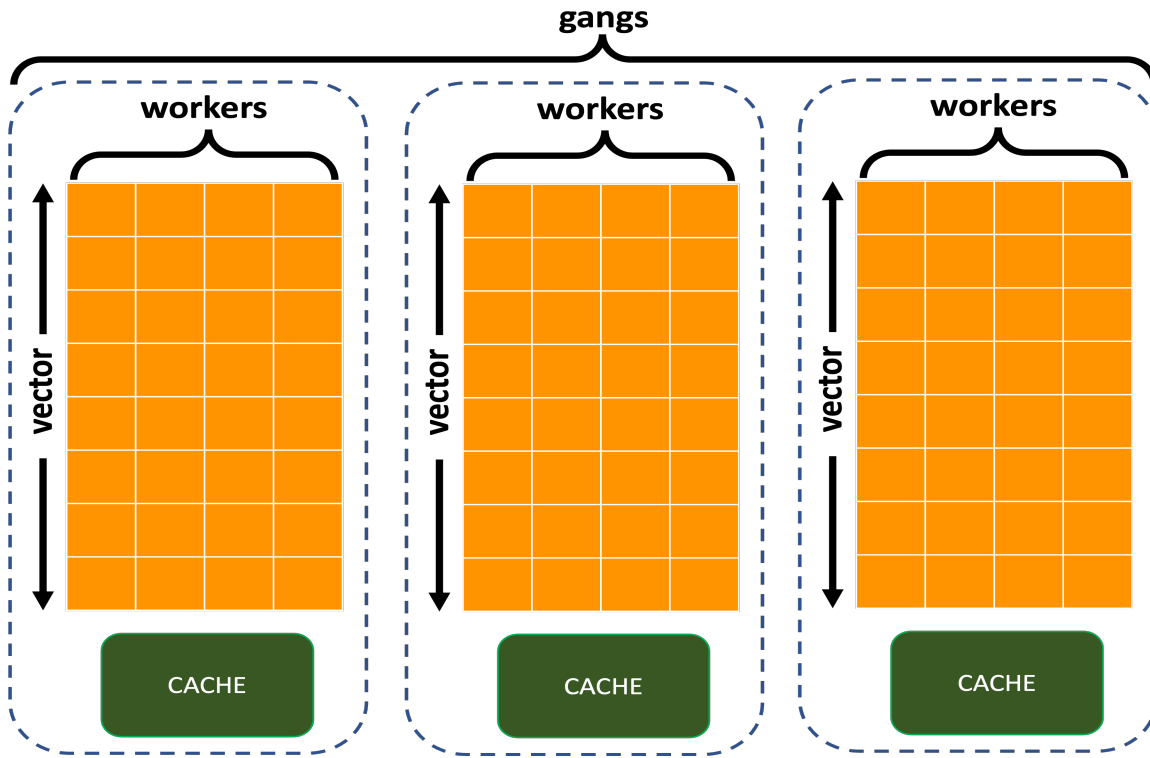


Figure 2. The map of gangs, workers, and vectors (adapted from Jiang et al., 2019)

```

1  MODULE W3WAVEMD
2      USE CONSTANTS
3      USE W3GDATMD
4      !! use other modules !!
5      REAL      :: Variables
6
7      !! ===== !!
8      !! BUNCH OF CODES !!
9      !! ===== !!
10
11     DO IT=ITO, NT !Time loop
12         !! ===== !!
13         !! BUNCH OF CODES !!
14         !! ===== !!
15
16         DO JSEA=1, NSEAL !grids loop
17             !! BUNCH OF CODES !!
18             CALL W3SRCEMD(variable list)
19             !! BUNCH OF CODES !!
20         END DO
21
22     !! ===== !!
23     !! BUNCH OF CODES !!
24     !! ===== !!
25     END DO !Time Loop
26     !! ===== !!
27     !! BUNCH OF CODES !!
28     !! ===== !!
29 END MODULE W3WAVEMD

```

(a)

```

1  MODULE W3WAVEMD
2      USE CONSTANTS
3      USE W3GDATMD
4      !! use other modules !!
5      REAL      :: Variables
6      !$acc declare create(needed variables on GPU)
7      !! BUNCH OF CODES !!
8      !$acc enter data copyin(variables to GPU)
9      DO IT=ITO, NT !Time loop
10         !! BUNCH OF CODES !!
11     !$acc update device(GPU data)
12     !$acc parallel copy(NSEAL) num_gangs(NSEAL)
13     !$acc loop gang
14         DO JSEA=1, NSEAL !grids loop
15             !! BUNCH OF CODES !!
16             CALL W3SRCEMD(variable list)
17             !! BUNCH OF CODES !!
18         END DO
19
20     !$acc end parallel
21     !$acc update host(needed variables on CPU)
22     !! ===== !!
23     !! BUNCH OF CODES !!
24     !! ===== !!
25     END DO !Time Loop
26     !! BUNCH OF CODES !!
27
28     !$acc exit data delete(all GPU variables)
29 END MODULE W3WAVEMD

```

(b)

Figure 3. A schematic representation of WAVEWATCH III Original FORTRAN source code for the W3WAVEMD module (a) and its OpenACC directives version (b)

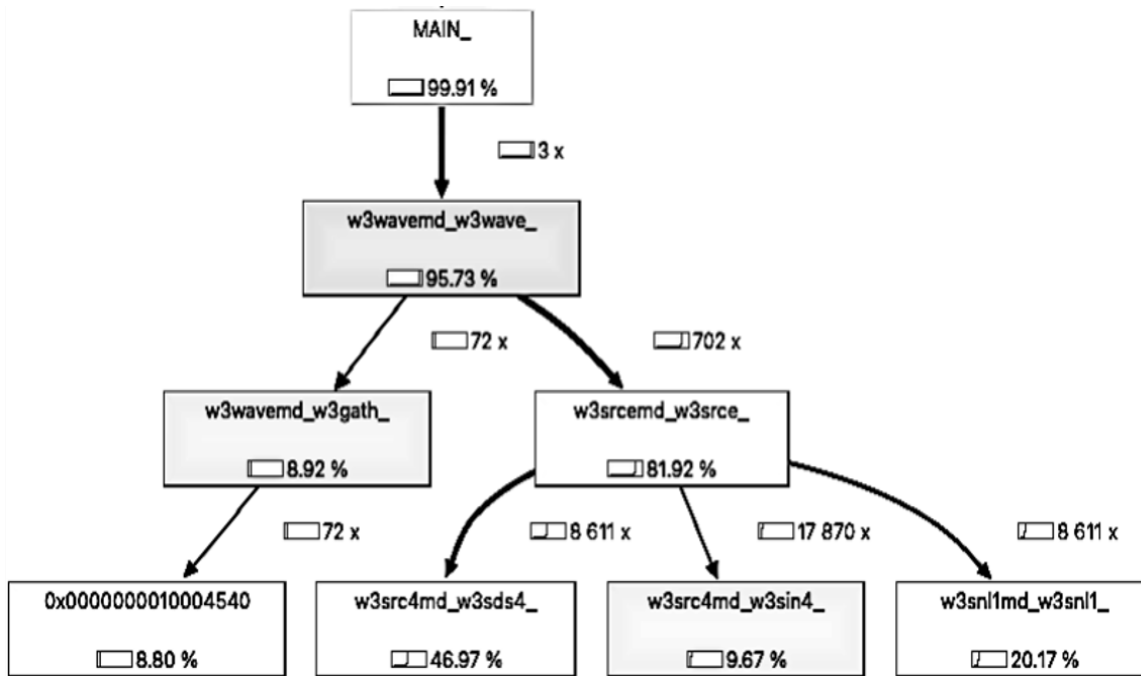


Figure 4. The Callgraph obtained by profiling WAVEWATCHIII with 300 CPU cores on Summit. Each box includes the name of the subroutine and its relative execution time as a percentage

Table 1. GPU hardware specifications.

	SUMMIT (V100)	KODIAK (P100)
Compute Capability	7	6
Global memory size	16 GB	16 GB
L1 cache	10 MB	1.3 MB
L2 cache	6 MB	4 MB
Shared memory size / SM	Configurable up to 96 KB	49 KB
Constant memory	64 KB	64 KB
Register File Size	256 KB (per SM)	256 KB (per SM)
32-bit Registers	65536 (per SM)	65536 (per SM)
Max registers per thread	255	255
Number of multiprocessors (SMs)	80	56
Warp size	32 threads	32 threads
Maximum resident warps per SM	64	64
Maximum resident blocks per SM	32	32
Maximum resident threads per SM	2048	2048
Maximum threads per block	1024	1024

Table 2. The speedups and simulation times associated with offloading one MPI rank exclusively to one GPU on Kodiak and Summit with 59K and 228K mesh sizes for MPI ranks ranging from 1 to 32

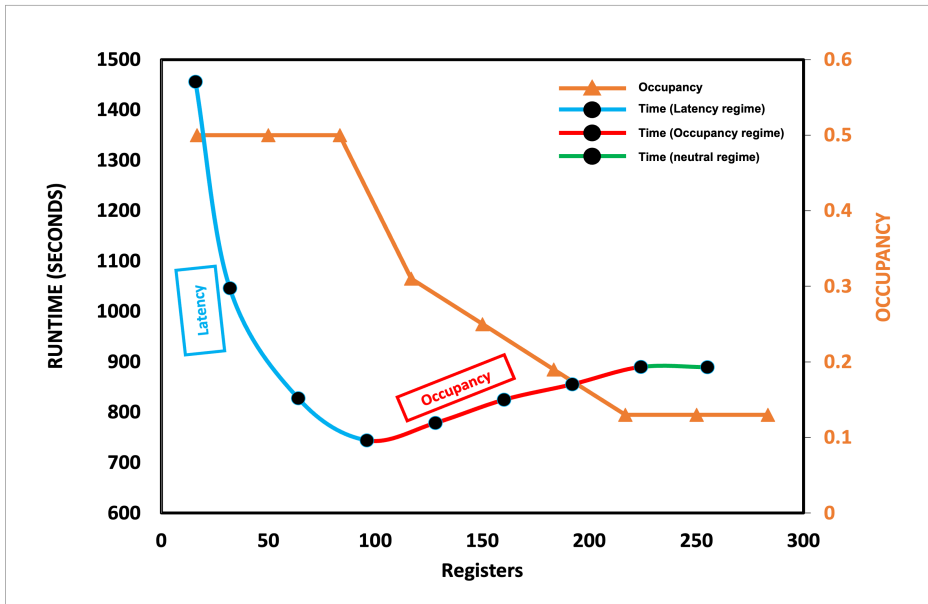
		MPI Ranks	CPU	GPU	Speedup
SUMMIT	59K	1	6688.97	1429.92	4.7
		2	3347.84	748.39	4.5
		4	1685.4	386.4	4.4
		8	842.91	205.45	4.1
		16	435.5	113.63	3.8
		32	236.84	64.39	3.7
	228K	1	38982.1	5905.57	6.6
		2	19827.3	3042.59	6.5
		4	9901.16	1544.22	6.4
		8	3424	778.45	4.4
		16	1720.87	399.59	4.3
		32	878.05	219	4.0
KODIAK	59K	1	5574.38	2477.38	2.3
		2	2856.38	1291.53	2.2
		4	1564.3	700.68	2.2
		8	873.05	373.78	2.3
		16	398.54	195.99	2.0
		32	239.64	114.2	2.1
	228K	1	22864.27	9636.75	2.4
		2	12184.45	5592.38	2.2
		4	6313.02	2964.1	2.1
		8	3423.55	1448.44	2.4
		16	1586.75	746.27	2.1
		32	871.22	381.01	2.3

Table 3. Node baseline comparison runtime (seconds) with different MPI ranks per GPU configuration on a Summit node. The bold numbers in [] are the speedup relative to the whole CPU on a summit node (42 MPI ranks). In the last row, we present the speedups and simulation times when comparing only the GPU-accelerated subroutine, W3SRCEMD.

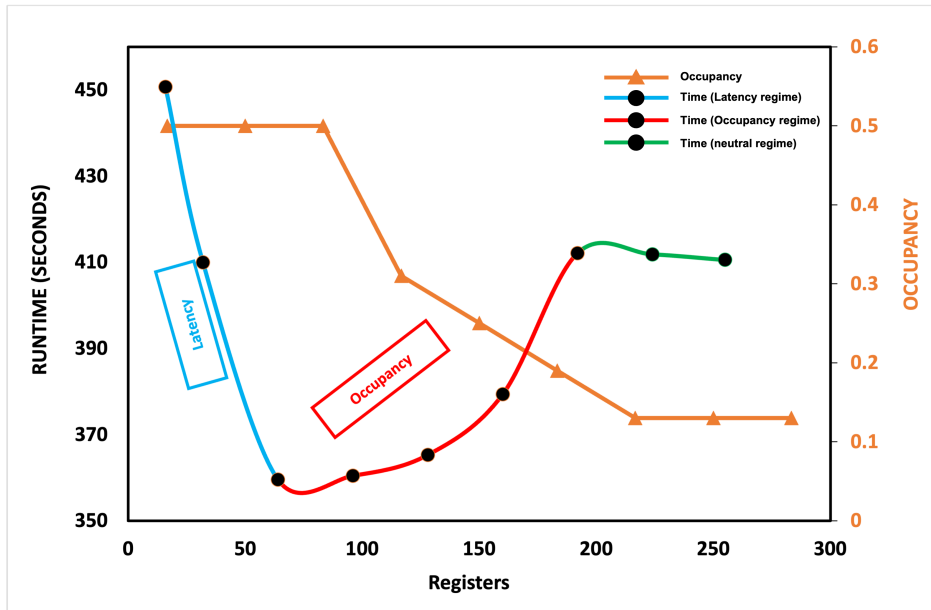
MESH SIZE	CPU	GPU			
	42 ranks	42 ranks (7/GPU)	24 ranks (4/GPU)	18 ranks (3/GPU)	12 ranks (2/GPU)
59K	180.72	133.18 [1.36x]	124.28 [1.45x]	127.17 [1.42x]	146.34 [1.23x]
228K	683.66	483.22 [1.41x]	523.91 [1.30x]	533.08 [1.24x]	612.60 [1.12x]
228K (W3SRCEMD)	589.56	389.12 [1.52x]	407.54 [1.45x]	366.00 [1.61x]	395.31 [1.49x]

Table 4. The speedup of multiple nodes for the 228k mesh size on summit. We used all the 42 CPU cores and 6 GPU on each node with a configuration of 7 MPI ranks per GPU.

NODES	Gridpoints (Per rank)	CPU	GPU	Speedup
1	5441	683.66	483.22	1.41x
2	2721	354.12	224.69	1.58x
3	1814	238.02	171.15	1.39x
4	1360	195.59	141.44	1.38x
5	1088	154.17	126.84	1.22x



(a)



(b)

Figure 5. WW3-W3SRCMD.gpu runtime (primary vertical axis) and occupancy (secondary vertical axis in orange) based on register counts on Kodiak (a) and Summit (b). For the runtime, the light blue part of the line represents the latency-dominant region, the red part represents the occupancy-dominant region, and the green part represents the neutral region. We analyzed the result of running 228K mesh size on 16 GPUs and CPUs.

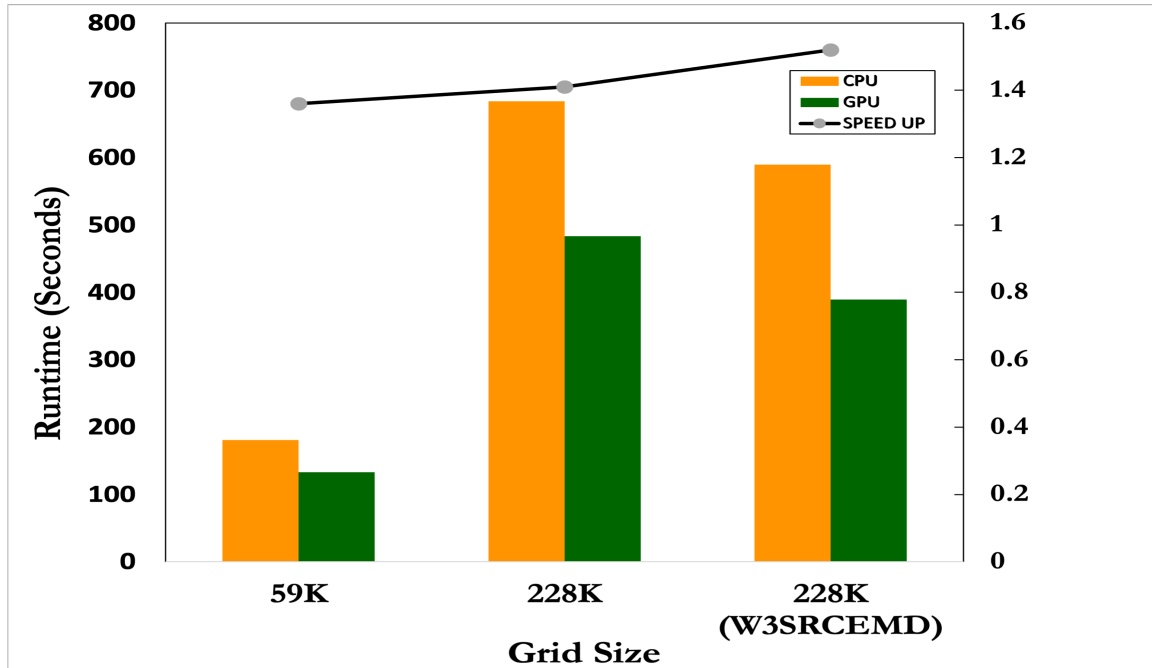


Figure 6. The runtime of WW3.cpu (orange) and WW3-W3SRCEMD.gpu (green) for 42 CPU cores and 7 MPI ranks per GPU on a Summit node with 59K and 228K meshes. The last bar chart is the runtimes for only the GPU-accelerated subroutine, W3SRCEMD on 228K grid size. The black line represents the speedups of WW3-W3SRCEMD.gpu over WW3.cpu.

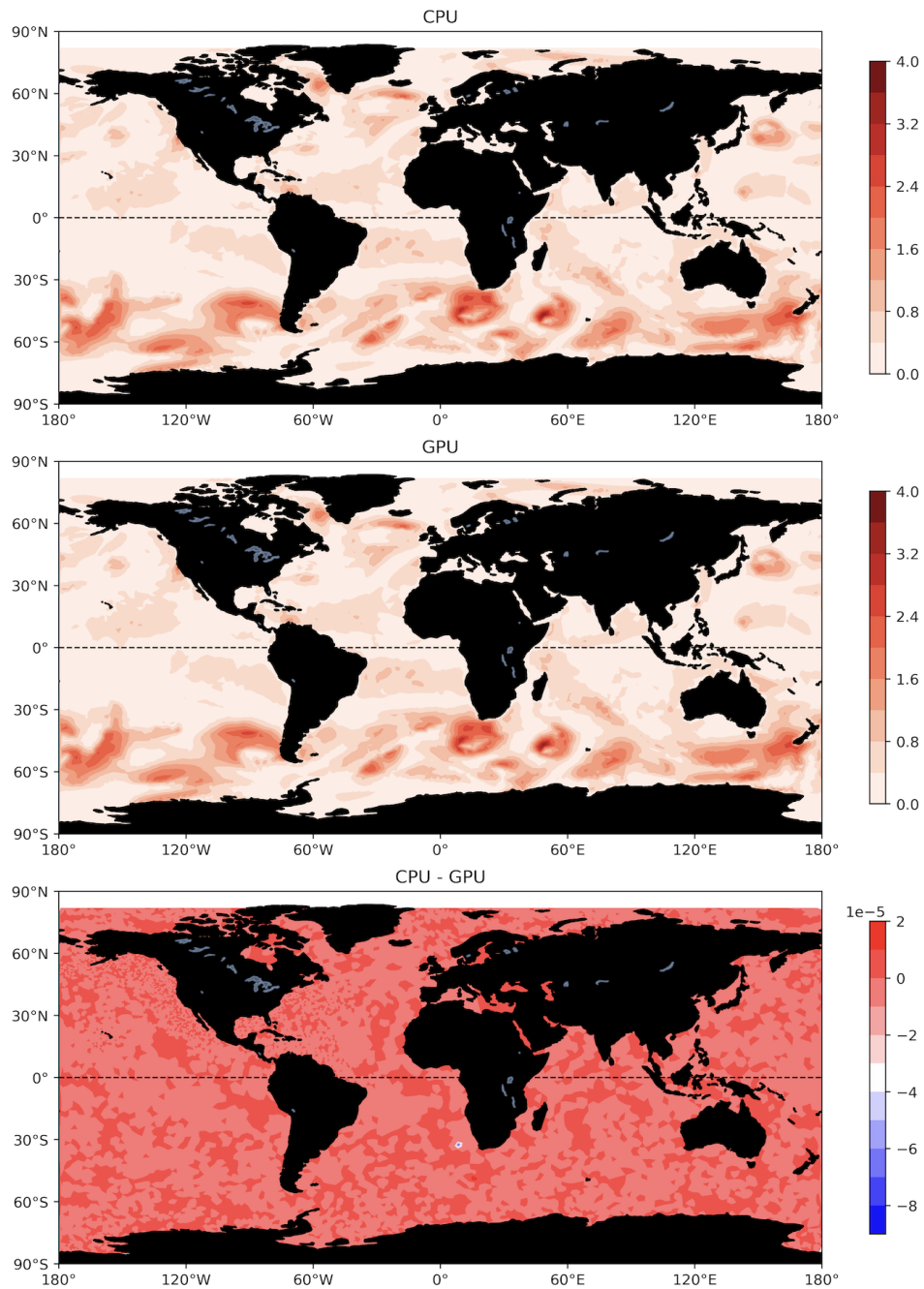


Figure 7. The average of WW3.cpu (top), WW3-W3SRCEMD.gpu (middle) last 2-hour simulations and their differences (bottom) for significant wave height

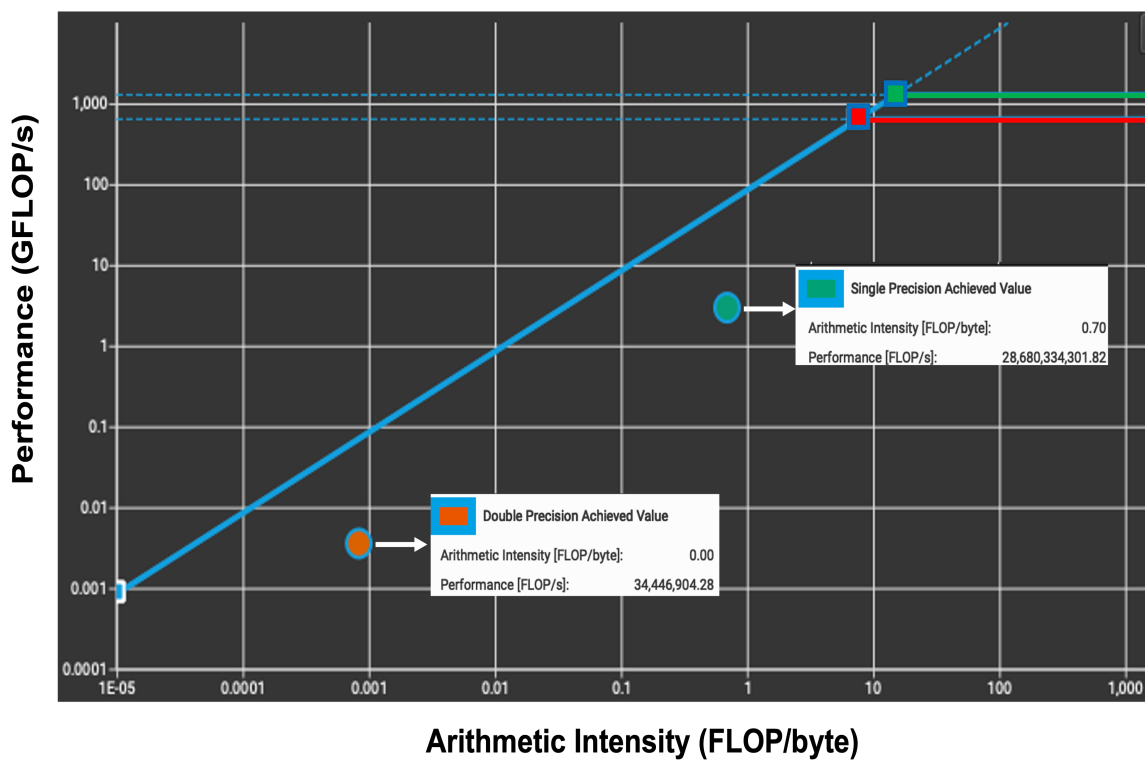


Figure 8. Roofline model for W3SRCEMD kernel from NVIDIA Nsight Compute. The red and green dots are the double and single precision values respectively.