# Comparing the Performance of Julia on CPUs versus GPUs and Julia-MPI versus Fortran-MPI: a case study with MPAS-Ocean (Version 7.1)

Robert R. Strauss[1], Siddhartha Bishnu[2], and Mark R. Petersen[2]

[1]Center for Nonlinear Studies, Los Alamos National Laboratory, NM, 87545, USA
[2]Computational Physics and Methods Group, Los Alamos National Laboratory, NM, 87545, USA

**Correspondence:** Mark R. Petersen (mpetersen@lanl.gov)

**Abstract.** Some programming languages are easy to develop at the cost of slow execution, while others are fast at run time but much more difficult to write. Julia is a programming language that aims to be the best of both worlds—a development and production language at the same time. To test Julia's utility in scientific high-performance computing (HPC), we built an unstructured-mesh shallow water model in Julia and compared it against an established Fortran-MPI ocean model, MPAS-Ocean,

5   as well as a Python shallow water code. Three versions of the Julia shallow water code were created, for: single-core CPU; graphics processing unit (GPU); and Message Passing Interface (MPI) CPU clusters. Comparing identical simulations revealed that our first version of the  Julia model was 13 times faster than Python using Numpy, where both used an unthreaded single-core CPU. Further Julia optimizations, including static typing and removing implicit memory allocations, provided an additional 10–20x speed-up of the single-core CPU Julia model. The GPU-accelerated Julia code is attained a speed-up

10  of 230-380x compared to the single-core CPU Julia code. Parallelized Julia-MPI performance was identical to Fortran-MPI MPAS-Ocean for low processor counts, and ranges from 2x faster to 2x slower for higher processor counts. Our experience is that Julia development is fast and convenient for prototyping, but that Julia requires further investment and expertise to be competitive with compiled codes. We provide advice on Julia code optimization for HPC systems.

## 1   Introduction

15  A major concern in computer modeling is the trade-off between execution speed and code development time. In general, programs in scripting languages like Python and Matlab are faster to develop due to their simpler syntax and more relaxed typing requirements, but are limited by slower execution time. On the other end of the spectrum, we have compiled languages like C/C++ and Fortran, which have been extensively used in scientific computing for many decades. Programs in such languages are blessed with faster execution time, but are cursed with stricter and more cumbersome syntax, leading to slower

20  development time. The Julia language strikes a balance between these two categories (Perkel, 2019). It is a compiled language with execution speed similar to C/C++ or Fortran, if carefully written with strict syntax (Lin and McIntosh-Smith, 2021; Gevorkyan et al., 2019). It is also equipped with a more convenient syntax and features, such as dynamic typing, to accelerate code development in prototyping. To this day, the majority of scientific computing models are programmed in compiled

languages, which execute fast but can take can take years to develop—for example, the first version of MPAS-Ocean required
25   three years (Ringler et al., 2013). In this paper, we investigate the feasibility of writing Julia codes for computational physics
simulations, since a Julia program can not only ensure high performance but also less development time in the initial stages.
We develop a shallow water solver in Julia and compare its performance to an equivalent Fortran code.

An additional complication in choosing the best language is that layers of libraries have been added to C/C++ and Fortran
to accommodate evolving computer architectures. For the past 25 years, computational physics codes have largely used the
30   Message Passing Interface (MPI) to communicate between CPUs on separate nodes that do not share memory, and OpenMP
to parallelize within a node using shared-memory threads. With the advent of heterogeneous nodes containing both CPUs and
GPUs, scientific programmers have several new choices: writing kernels directly for GPUs in CUDA (Bleichrodt et al., 2012;
Zhao et al., 2017; Xu et al., 2015); adding OpenACC pragmas for the GPUs (Jiang et al., 2019); or calling libraries such as
Kokkos (Trott et al., 2022) that execute code optimized for specialized architectures on the back-end, while providing a simpler
35   front-end interface for the domain scientist. All of these require additional expertise, and add to the length and complexity of
the code base. Julia also provides an MPI library for parallelization across nodes in a cluster, and a CUDA library to parallelize
over GPUs within a node. We have written shallow water codes in Julia that adopt each of these parallelization strategies.

In recent years, shallow water solvers such as Oceananigans.jl (Ramadhan et al., 2020) and ShallowWaters.jl (Klöwer
et al., 2022) have been developed in Julia. These codes employ structured rectilinear meshes to discretize their domain,
40   and are equipped with capabilities for running on GPUs to achieve high performance. Here we conduct a comparison on
unstructured-mesh models, using the Fortran code MPAS-Ocean (Ringler et al., 2013) as a point of reference. MPAS-Ocean
employs unstructured near-hexagonal meshes with variable resolution capability and is parallelized with MPI for running on
supercomputer clusters. We developed a Julia model employing the same spatial discretization of MPAS-Ocean, and capable of
running in serial mode on a single core, or in parallel mode on a supercomputer cluster or a graphics card. We discuss the subtle
45   details of our implementations, compare the speed-ups attained, and describe the strategies employed to enhance performance.

## 2   Methods

### 2.1   Equation Set & TRiSK-Based Spatial Discretization

Our Julia model solves the shallow water equations (Cushman-Roisin and Beckers, 2011) in vector-invariant form. This is
sufficiently close to the governing equations for ocean and atmospheric models to be used as a proxy to test performance with
50   new codes and architectures. The equation set is

$$\boldsymbol{u}_t + qh\boldsymbol{u}^{\perp} = -g\nabla\eta - \nabla K, \tag{1a}$$

$$\eta_t + \nabla \cdot (h\boldsymbol{u}) = 0, \tag{1b}$$

where $\boldsymbol{u}$ is the horizontal velocity vector, $\boldsymbol{u}^{\perp} = \boldsymbol{k} \times \boldsymbol{u}$, $h$ is the layer thickness, $\eta$ is the surface elevation or sea surface height
(SSH), $K = |\boldsymbol{u}|^2/2$ is the kinetic energy, and $g$ is the acceleration due to gravity. If $b$ represents the topographic height and $H$
55   the mean depth, then $\eta = h + b - H$. Moreover, if $f$ denotes the Coriolis parameter, and $\zeta = \boldsymbol{k} \cdot \nabla \times \boldsymbol{u}$ the relative vorticity, then

the absolute vorticity, $\omega_a = \zeta + f$, and the potential vorticity, $q = \omega_a/h$. The term $qh\boldsymbol{u}^\perp$ is the thickness flux of the PV in the direction perpendicular to the velocity, rotated counterclockwise on the horizontal plane. Ringler et al. (2010) refer to it as the non-linear Coriolis force since it consists of the quasi-linear Coriolis force $f\boldsymbol{u}^\perp$ and the rotational part $\zeta\boldsymbol{u}^\perp$ of the non-linear advection term $\boldsymbol{u}\cdot\nabla\boldsymbol{u}$. We spatially discretize the prognostic equations in (1) using a mimetic finite volume method based on the

60  TRiSK scheme, which was first proposed by (Thuburn et al., 2009), and then generalized by (Ringler et al., 2010). This method was chosen to horizontally discretize the primitive equations of MPAS-Ocean while invoking the hydrostatic, incompressible, and Boussinesq approximations on a staggered C-grid. Since this horizontal discretization guarantees conservation of mass, potential vorticity, and energy, it makes MPAS-Ocean a suitable candidate to simulate mesoscale eddies.

Our spatial domain is tessellated by two meshes, a regular planar hexagonal primal mesh and a regular triangular dual mesh.

65  Each corner of the primal mesh cell coincides with a vertex of the dual mesh cell and vice versa. A line segment connecting two primal mesh cell centers is the perpendicular bisector of a line segment connecting two dual mesh cell centers and vice versa. Regarding our prognostic variables, the scalar SSH $\eta$ is defined at the primal cell centers, and the normal velocity vector $\boldsymbol{u}_e$ is defined at the primal cell edges. The divergence of a two-dimensional vector quantity is defined at the positions of $\eta$, while the two-dimensional gradient of a scalar quantity is defined at the positions of $\boldsymbol{u}_e$ and oriented along its direction. The

70  curl of a vector quantity is defined at the vertices of the primal cells. Finally, the tangential velocity $\boldsymbol{u}_e^\perp$ along a primal cell edge is computed diagnostically using a flux mapping operator from the primal to the dual mesh, which essentially takes a weighted average of the normal velocities on the edges of the cells sharing that edge. Interested readers may refer to Thuburn et al. (2009) and Ringler et al. (2010) for further details on the mesh specifications.

At each edge location $\boldsymbol{x}_e$, two unit vectors $\boldsymbol{n}_e$ and $\boldsymbol{t}_e$ are defined parallel to the line connecting the primal mesh cells, and in

75  the perpendicular direction rotated counterclockwise on the horizontal plane, such that $\boldsymbol{t}_e = \boldsymbol{k} \times \boldsymbol{n}_e$. The discrete equivalent of the set of equations (1) is

$$(u_e)_t = F_e^\perp \widehat{q}_e - \left[\nabla(g\eta)_i + K_i\right]_e, \tag{2a}$$

$$(\eta_i)_t = -\left[\nabla \cdot F_e\right]_i, \tag{2b}$$

where $F_e = \widehat{h_e}u_e$ and $F_e^\perp$ represent the thickness fluxes per unit length in the $\boldsymbol{n}_e$ and $\boldsymbol{t}_e$ directions respectively. The layer

80  thickness $h_i$, the SSH $\eta_i$, the topographic height $b_i$, and the kinetic energy $K_i$ are defined at the centers $\boldsymbol{x}_i$ of the primary mesh cells, while the velocity $u_e$ are defined at the edge points $\boldsymbol{x}_e$. The symbol $\widehat{(.)}_e$ represents an averaging of a field from its native location to $\boldsymbol{x}_e$. The discrete momentum equation (2b) is obtained by taking the dot product of (1b) with $\boldsymbol{n}_e$, which modifies the non-linear Coriolis term to

$$\boldsymbol{n}_e \cdot \widehat{q}_e\widehat{h_e}\boldsymbol{u}^\perp = \widehat{q}_e\widehat{h_e}\boldsymbol{n}_e \cdot (\boldsymbol{k} \times \boldsymbol{u}) = \widehat{q}_e\widehat{h_e}\boldsymbol{u} \cdot (\boldsymbol{n}_e \times \boldsymbol{k})$$

85  $$= -\widehat{q}_e\widehat{h_e}\boldsymbol{u} \cdot \boldsymbol{t}_e = -\widehat{q}_e\widehat{h_e}u_e^\perp = -F_e^\perp \widehat{q}_e. \tag{3}$$

Given the numerical solution at time level $t^n = n\Delta t$, with $\Delta t$ representing the time step and $n \in \mathbb{Z}_{\geq 0}$, the Julia model first computes the time derivative or tendency terms of (2) as functions of the discrete spatial and flux-mapping operators of the

TRiSK scheme. Then it advances the numerical solution to time level $t^{n+1}$ using the forward-backward method

$$u^{n+1} = u^n + \Delta t \mathcal{F}(u^n, h^n), \tag{4}$$

$$h^{n+1} = h^n + \Delta t \mathcal{G}\left(u^{n+1}, h^n\right), \tag{5}$$

where $\mathcal{F}$ and $\mathcal{G}$ represent the discrete tendencies of the normal velocity and the layer thickness in functional form, and the subscripts representing the positions of these variables have been dropped for notational simplicity.

The following sections introduce the new codes that were created for this study. Three versions of the Julia code were written (Strauss, 2023): the base single-core CPU version, an altered version for GPUs with CUDA, and a multi-node CPU implementation with Julia-MPI. These were compared against existing Fortran-MPI and single-core Python versions of shallow-water TRiSK models. All use a standard MPAS unstructured-mesh file format that specifies the geometry and topology of the mesh, and includes index variables that relate neighboring cells, edges, and vertices. All models have an inner (fastest-moving) index for the vertical coordinate and were tested with 100 vertical layers to mimic performance in a realistic ocean model.

### 2.2 Single-Core CPU Julia Implementation

The serial-mode implementation on a single core involves looping over every cell and edge of the mesh to (a) compute the tendencies, i.e. the right-hand side terms of the prognostic equations (2) and (b) advance their values to the next time step. The tendencies can be functions of the dependent and independent variables as well as spatial derivatives of the dependent variable. The serial version of our model is the simplest one from the perspective of transforming the numerical algorithms into code.

In order to highlight differences in formulation, we provide a Julia code example for the single tendency term from (2) for the SSH gradient $-g\nabla\eta$, which is discretized as $-\left[g\nabla\eta_i\right]_e$. We then add a vertical index $k$ to mimic the performance of a multi-layer ocean model, but each layer is trivially redundant. In a full ocean model this term would be the pressure gradient, and would involve the computation of pressure as a function of depth and density. For the single-core CPU version, the Julia function computing the SSH gradient is

**Listing 1.** Julia example for serial CPU

```
1: velocity_tendencies!(sshGradient, ssh, ...)
2:
3: function velocity_tendencies!(sshGradient, ssh, ...)
4:     for iEdge in 1:nEdges
5:         cell1 = cellsOnEdge[1,iEdge]
6:         cell2 = cellsOnEdge[2,iEdge]
7:         for k in 1:nVertLevels
8:             sshGradient[k,iEdge] = - gravity / dcEdge[iEdge]
9:                 * ( ssh[k,cell2] - ssh[k,cell1] )
10:        end
11:    end
```

120    12: **end**

Here `cellsOnEdge` is an array of index variables describing the mesh that points to the cells neighboring an edge, and `dcEdge` represents the distance between the centers of adjacent cells sharing the edge on which the normal velocity tendency is computed. In the actual code all the tendency terms are computed within this function, but here we only show the ssh gradient as a brief sample.

125    ## 2.3  SIMD GPU Julia Implementation

GPUs are very powerful tools for SIMD (Same Instruction Multiple Data) computations: they have thousands of independent threads, which can execute the same operation at the same time with different input values. Since we numerically solve the same prognostic equation for (a) the SSH at every cell center $x_i$, and (b) the normal velocity at every edge $x_e$ of the mesh, a GPU is a logical tool to employ for our computations. By placing subsets of cells and edges on different threads of the GPU,

130    we can perform the tendency computations, and advance the prognostic variables at once in parallel rather than looping over every cell and edge, which would scale in wall-clock time according to the size of the mesh.

We wrote CUDA kernels for an Nvidia GPU using the CUDA.jl library for computing the tendencies and advancing the prognostic variables to the next time step. The code for the single-core implementation can be converted to CUDA with surprising ease by removing the `for` loop over the cells and edges, and instead performing the underlying computation on a

135    single cell or edge:

**Listing 2.** Julia example for GPU with CUDA

```
1: CUDA.@cuda blocks=cld(nEdges, 1024) threads=1024 maxregs=64
2:    velocity_tendencies_cuda!(sshGradient, ssh, ...)
3:
4: function velocity_tendencies_cuda!(sshGradient, ssh, ...)
5:    iEdge = (CUDA.blockIdx().x - 1) * CUDA.blockDim().x
6:        + CUDA.threadIdx().x
7:    cell1 = cellsOnEdge[1,iEdge]
8:    cell2 = cellsOnEdge[2,iEdge]
9:    for k in 1:nVertLevels
10:       sshGradient[k,iEdge] = - gravity / dcEdge[iEdge]
11:           * ( ssh[k,cell2] - ssh[k,cell1] )
12:    end
13: end
```

Each cell and edge of the mesh will be designated to a different thread on the GPU. The computation for a single cell or edge

150    will run on a single thread, and a CUDA method will be used to map the index of the thread to the index of the cell ($i$) or edge ($e$), at which the prognostic variable is being updated. To execute this method over all threads on the GPU, we use a CUDA

macro to call our kernel and designate the number of threads to use, which should be equal to the number of cells or edges in the mesh. Note that the inner computation of `pressureGradient` is identical for the CPU and CUDA kernal codes.

## 2.4 CPU/MPI Julia Implementation

155   Rather than iterating through every cell or edge of the mesh, we may parallelize the simulation with multiple processors by assigning to each processor a portion of the mesh, a process called domain decomposition. However, the computations of some spatial operators may require information from the outermost cells of the adjacent processors. So, we need the neighboring processors to communicate these pieces of information with each other. To ensure an efficient communication, we include an extra ring or "halo" of cells around the boundary of the region assigned to each processor, which overlaps with the

160   region assigned to adjacent processors. We do not compute the tendencies of the prognostic variables in the halo region of a processor. In fact, we cannot perform this operation without information in an additional ring of halo cells, which is not assigned to the processor under consideration. So, we obtain the updated values of the prognostic variables in the halo region by communication with adjacent processors, which contain these halo cells in their interior, and update the prognostic variables in them.

165   A number of crucial modifications are necessary to implement this parallelization scheme. For instance, the simulation methods are amended so that each process (rank) only performs computations for the set of cells or edges assigned to it. We use the MPI communication channel (comm) to receive the updated values of the prognostic variables in the halo region of a processor from adjacent processors which advance these variables. Similarly, we send the updated values of the prognostic variables along the outermost region of the above-mentioned processor to adjacent processors, for which these variables belong

170   in the halo regions. For the TRiSK-based spatial discretization and the forward-backward time-stepping method, the halo region consists of only one layer (one halo ring) of cells.

**Listing 3.** Julia example for CPU with MPI

```
1: # each process executes the following, receiving a different value
2: # on each rank:
3: comm = MPI.COMM_WORLD
4: rank = MPI.Comm_rank(comm)
5:
6: myCells = cells_for_rank(mesh_file, rank, partition_file)
7: myEdges, myHaloEdges = edges_on_cells(myCells)
8:
9: velocity_tendencies!(myEdges, sshGradient, ssh, ...)
10: update_halo_edges!(sshGradient, myHalodEdges, rank, comm)
11:
12: function velocity_tendencies!(myEdges, sshGradient, ssh, ...)
13:     for iEdge in myEdges
```

```
185  14:         cell1 = cellsOnEdge[1,iEdge]
     15:         cell2 = cellsOnEdge[2,iEdge]
     16:         for k in 1:nVertLevels
     17:             sshGradient[k,iEdge] = - gravity / dcEdge[iEdge]
     18:                 * ( ssh[k,cell2] - ssh[k,cell1] )
190  19:         end
     20:     end
     21: end
     22:
     23: function update_halo_edges!(data, edgesInMyHalo, rank, comm)
195  24:     for neighborRank in find_neighbors(rank, comm)
     25:         MPI.Irecv!(data[edgesInMyHalo,:], neighborRank, 0, comm)
     26:         edgesToNeighbor = find_halo_overlap(rank, neighbor, comm)
     27:         MPI.Isend(data[edgesToNeighbor,:], neighborRank, 0, comm)
     28:     end
200  29: end
```

Here `myCells` and `myEdges` are the lists of cells and edges in the local domain, owned by the rank running this code, plus its halo.

## 2.5 CPU/MPI Fortran Implementation

j The baseline comparison code for this study is the Model for Prediction Across Scales (MPAS-Ocean) (Ringler et al., 2013; Petersen et al., 2015), which is written in Fortran with MPI communication commands. It is the ocean component of the Energy Exascale Earth System Model (E3SM) (Golaz et al., 2019; Petersen et al., 2019), the climate model developed by the US Department of Energy. In this study, the code is reduced from a full ocean model solving the primitive equations to simply solving for velocity and thickness (1). Thus the majority of the code is disabled, including the tracer equation, vertical advection and diffusion, the equation of state, and all parameterizations. In order to match the Julia simulations, we employ a forward-backward time-stepping scheme, exchange one-cell-wide halos after each time step, compute 100 layers in the vertical array dimension, and use the identical Cartesian hexagon-mesh domains (Petersen et al., 2022).

MPAS-Ocean is an excellent comparison case for Julia because it is a well-developed code base that uses Fortran and MPI, which have been standard for computational physics codes since the late 1990s. The highest resolution simulations in past studies used over three million horizontal mesh cells and 80 vertical layers, scale well to tens of thousands of processors (Ringler et al., 2013) and have been used for detailed climate simulations (Caldwell et al., 2019). MPAS-Ocean includes OpenMP for within-node memory access, and is currently adding OpenACC for GPU computations, but these were not used for this comparison to Julia-MPI on a CPU cluster.

## 2.6 Single-Core CPU Python Implementation

In addition to MPAS-Ocean, we compare the performance of the Julia shallow water code against an object-oriented single-core
Python code (Bishnu, 2022), which uses Numpy. The Python code solves the rotating shallow water system of equations
using two types of spatial discretizations: the TRiSK-based mimetic finite volume method used in MPAS-Ocean, and a
discontinuous Galerkin Spectral Element Method (DGSEM). The code offers a number of standard predictor-corrector and
multistep time-stepping methods, including those analyzed for ocean modeling in Shchepetkin and McWilliams (2005).

The Julia shallow water code was first written by translating this Python code into Julia syntax. While the Julia code was
expanded for parallelization and performance, the Python code was further developed to serve as a platform for conducting
a verification suite of shallow water test cases for the barotropic solver of ocean models. Each of these test cases in the
Python code verifies the implementation of a subset of terms in the prognostic momentum and continuity equations, e.g. the
linear pressure gradient term, the linear constant or variable-coefficient Coriolis and bathymetry terms, and the non-linear
advection terms. Bishnu et al. (2022) and Bishnu (2021) provide detailed discussions on these test cases along with specifics
of the numerical implementation, the time evolution of the numerical error for both spatial discretizations and a subset of the
time-stepping methods, and results of convergence studies with refinement in both space and time, only in space, and only in
time. Out of all of these test cases, only the linear coastal Kelvin wave and inertia-gravity wave test cases were implemented
in the Julia code for the current study.

While not used in this study, a number of libraries exist to accelerate Python for various architectures. These include Numba
and PyCuda for GPUs, mpi4py for CPU clusters, and Cython for single-CPU acceleration. Numba (Lam et al., 2015) is an
open-sourced Anaconda-sponsored NumPy-aware optimizing compiler, which translates Python functions to fast machine
code at runtime using the remarkable industry-standard LLVM compiler library. PyCUDA (Klöckner et al., 2012), written in
C++ (the base layer) and Python, provides access to Nvidia's CUDA parallel computation API from Python. Mpi4py (Dalcín
et al., 2005, 2008), provides Python bindings for the Message Passing Interface (MPI) standard. As an alternative, one can
'cythonize' an existing Python code by providing static type declarations and class attributes, that can then be translated
to C++/C code and to C-Extensions for Python. Cython is an optimising static compiler for both the Python programming
language and the extended Cython programming language. It is designed to offer C-like performance with code mostly written
in Python with additional C-inspired syntax. The rotating shallow water Python code Bishnu (2022) is currently undergoing
cythonization. Cythonized codes can further be accelerated on GPUs using Nvidia's HPC C++ compiler, and the C++ Standard
Parallelism (stdpar) for GPUs (Srinath). However, the extent of additional modifications and enhancements required to bring
GPU-accelerated C++ algorithms to the Python ecosystem may not always be a reasonable investment of time. As we will
see in later sections, a serial Julia code, which already achieves the performance of a fast compiled language, does not
require extensive modifications to be parallelized on GPUs or multiple cores, and is therefore more convenient than python for
high-performance scientific computing applications.

## 3 Results

### 3.1 Model Verification

Each serial and parallel implementation of the shallow water model described in the previous section was verified for accuracy with convergence tests against exact solutions. We obtained the expected second-order convergence of the various TRiSK-based spatial operators on a uniform planar hexagonal MPAS-Ocean mesh. The operators included the gradient, the divergence, the curl, and the flux-mapping operator used to interpolate the tangential velocities from the normal velocities (Figure 1). The formulation of these operators is shown in Figure 3 of Ringler et al. (2010). Once the operator tests were complete, the linearized shallow water equations were verified against exact solutions for the coastal Kelvin wave and inertia-gravity wave cases, as described in Bishnu et al. (2022) and Bishnu (2021). With refinement in both space and time, we observe the expected first-order convergence of the numerical solution (Figure 1), spatially discretized with the second-order TRiSK scheme, and advanced with the first-order forward-backward time-stepping method (Bishnu, 2021).

### 3.2 Acceleration of Julia with GPU Hardware

The Julia serial CPU version of the shallow water model was compared against the Julia CUDA library GPU version and the reference Python CPU code (Table 1 and Figure 2). Tests were conducted on the Darwin cluster at Los Alamos National Laboratory, using a single node equipped with Intel Cascade Lake CPUs (Gold 6248 with a clock rate of 2.5 GHz and 27.5M Cache) and the Nvidia Quadro RTX 8000 "Turing" GPU architecture (4608 CUDA cores, 16.3 TFLOPS peak single precision performance, 48 GB GPU memory, and GPU memory bandwidth of 672 GB/s). All performance tests described in this and the following sections used the coastal Kelvin wave test case on a planar hexagon mesh with the linear shallow water equations and 100 vertical layers. Samples are averaged over ten trials. All codes use double-precision (8 byte) real numbers, and performance tests do not include the time for initialization, input/output, or generating plots.

In our first version of the Julia single-core CPU code, we did not take any special steps for code optimization, and it was already 13 times faster than Python. Julia and Python both have dynamic typing, but Julia has the ability to go much faster since it also supports concrete typing. Julia is compiled, but hides it cleverly by compiling on the fly based on what datatypes are provided at run time. It supports a hierarchical abstract typing system, allowing for semi-specified types, such as "Any", which all types extend and is the default if no type is specified (thus acting like python), or "AbstractArray", which can be occupied at run time with any Array-like data.

After the initial Julia development, further effort was put into optimization, which led to a 10–20 times speed-up for the CPU-serial code. The changes included optimizing for memory management by tracking down and reducing unnecessary allocations that contributed significantly to the run time, as well as making all types and subtypes concrete rather than abstract, to minimize on-the-fly compilation. These improvements are explained in more detail in section 4.

We found the CUDA GPU implementation to be *significantly* faster than the single-core implementation. Because the memory transfer between the CPU and GPU takes many orders of magnitude longer than the actual on-GPU computations, we split them out in Table 1 and Figure 2. The memory transfers require between 0.015s and 0.68s and scale with the array
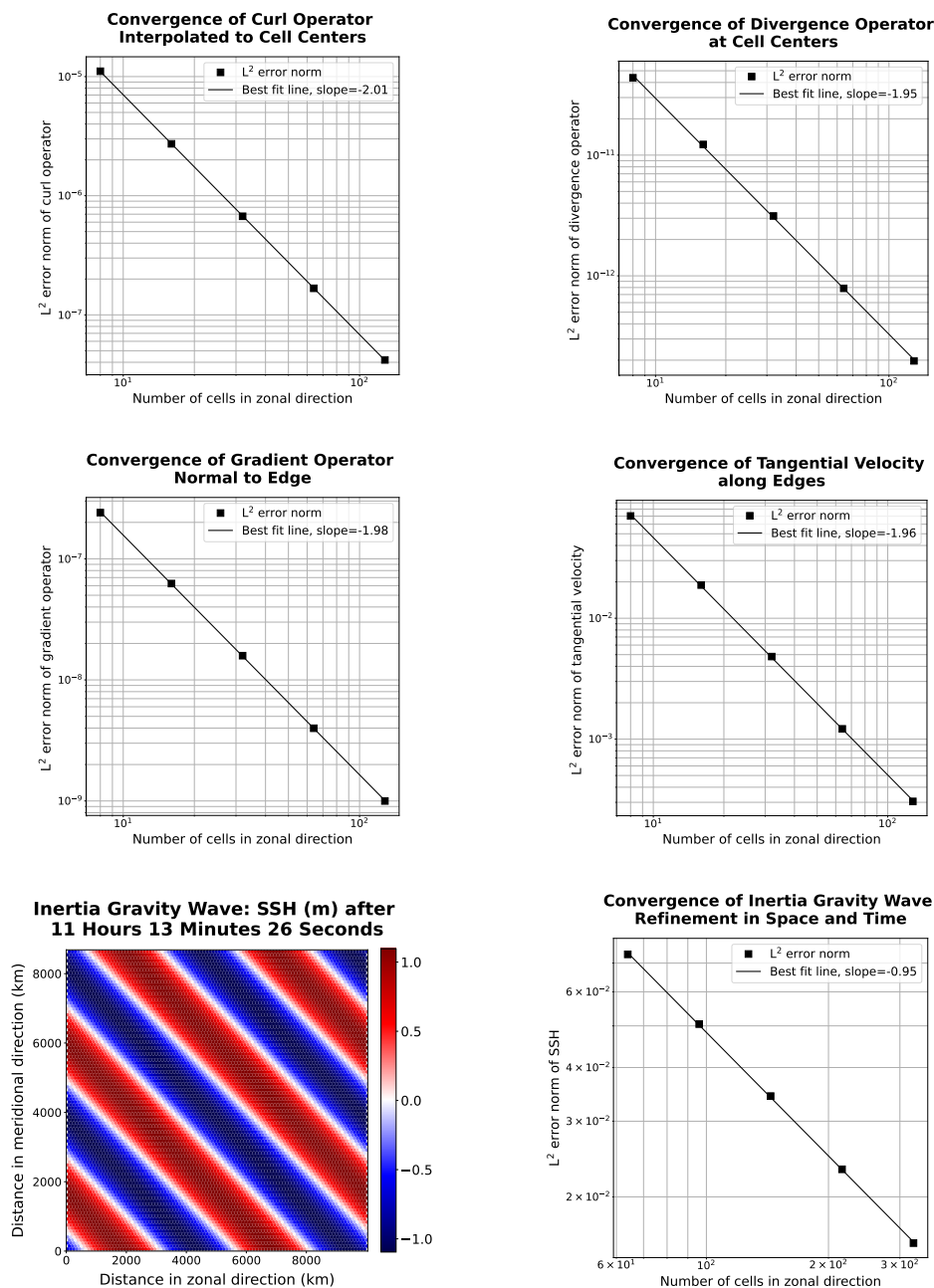
**Figure 1.** The first two rows show convergence plots of the TRiSK-based spatial operators for the newly-developed Julia code. Tests were run with both CPU and GPU implementations, and identical results were obtained. The slope of $-2$ indicates the expected second-order convergence. The third row shows a snapshot of the inertia-gravity wave test case, and the convergence plot of the numerical solution with refinement in both space and time.

size, while the GPU computations alone are 100 times faster, at 0.00027s for the 512x512 resolution case, and do not scale with resolution. This shows the power of GPUs, where computations alone can run over 40,000 times faster on the GPU than the CPU, but this speed-up is substantially diminished by the memory transfer time. Still, codes that are designed with a small memory footprint and limited memory transfer can greatly benefit from GPU computations. Strategically reducing array precision to 4-byte or even 2-byte reals for certain variables allows higher-resolution domains to fit on GPUs (Ye et al., 2022; Klöwer et al., 2022). In addition, single-precision floating point numbers (CUDA `Float32` data type) calculations may execute significantly faster than `Float64` (Julia Development Team, a). We did not leverage `Float32` in this work, but it shows that GPU simulations could run even faster than the results shown here.

Summing the GPU memory transfer and compute for the 10 timestep performance test, the GPUs were 229 to 386 times faster than the single CPU (Table 2). This compares to published studies of ocean models that show a speed-up from CPU to GPU ranging from 5–50 (Bleichrodt et al., 2012; Zhao et al., 2017; Xu et al., 2014), and a speed-up of up to 1556x for a GPU/CUDA Based Parallel Weather and Research Forecast Model (WRF) (Mielikainen et al., 2012). Note that our speed-up factor could be increased substantially by transferring data from the GPU to CPU less frequently. For a low-resolution ocean model with 30-minute time steps, the speed-ups in Table 2 correspond to collecting data every 10 time-steps, which is 5 hours of model time. One could instead collect data for analysis every 100 time-steps (~2 days), and that would result in a GPU speed-up of 2290 to 3860, because the compute time is negligible compared to the memory transfer. On the other hand, if model communication is required frequently for surface data forcing or coupling with atmospheric and sea ice components, the speed-up is drastically reduced. For example, if memory must be transferred between the CPU and GPU every time step, the speed-ups range from 23—39. The point is that GPU performance is wholly dependant on the GPU communication frequency.

| | 128x128 | 256x256 | 512x512 |
|---|---|---|---|
| Python, CPU | 3.08E+03 | 1.31E+04 | 4.96E+04 |
| Julia, CPU-serial (unoptimized) | 2.25E+02 | 8.64E+02 | 3.86E+03 |
| Julia, CPU-serial (optimized) | 1.12E+01 | 7.43E+01 | 3.33E+02 |
| Julia, GPU, total | 4.90E−02 | 2.03E−01 | 8.64E−01 |
| transfer to GPU | 2.98E−02 | 1.16E−01 | 4.58E−01 |
| compute on GPU | 2.51E−04 | 2.67E−04 | 2.67E−04 |
| transfer back to CPU | 1.53E−02 | 9.54E−02 | 6.84E−01 |

**Table 1.** Wall clock duration (seconds) of performing ten timesteps with 100 layers on an Intel Cascade Lake CPU or an NVidia Turing GPU.

GPU threads are grouped into threadblocks (or just "blocks") for efficiency. While calling the kernel function, we must specify the number of blocks and number of threads per block (the "block size"), as shown in listing 2. Within the kernel, we obtain the index of the block and thread, multiply the block index by the block size, and add the thread index to compute a global index. There is a maximum possible block size, but we can choose any smaller value to execute the kernel with. The
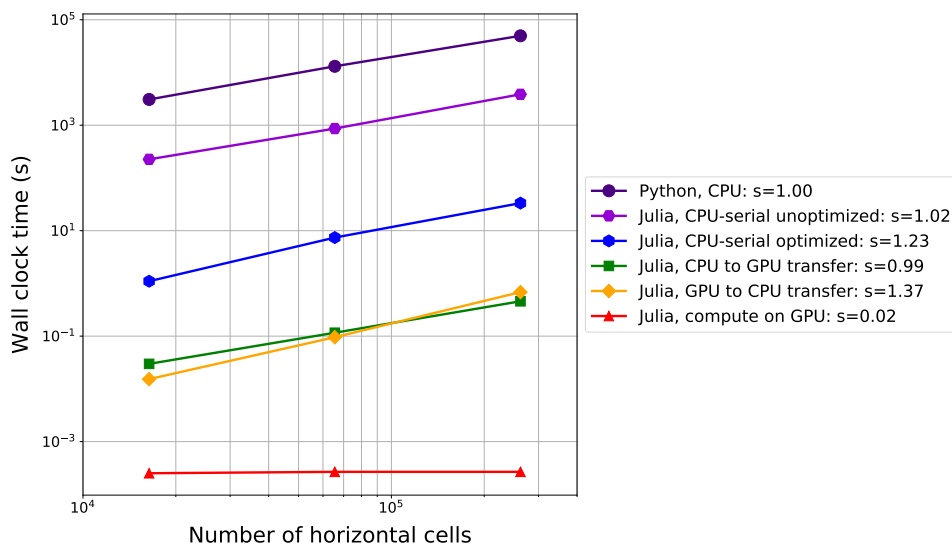
**Figure 2.** Timing data from Table 1, comparing ten timesteps of the Kelvin Wave test case on an Intel Cascade Lake CPU or an NVidia Turing GPU. The log-log slope, shown as s in the legend, is 1.0 for perfect scaling.
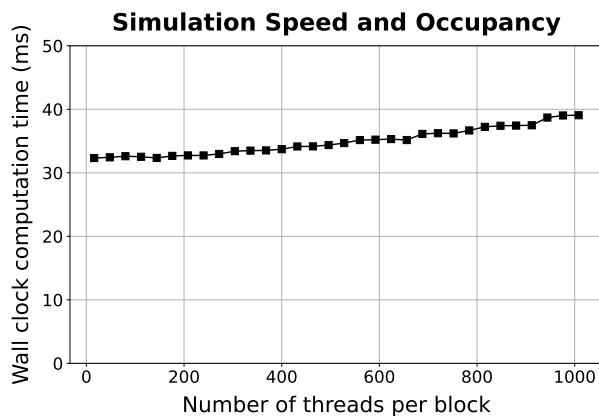


**Figure 3.** The same kernel was executed with the same data but different block sizes and the average execution time over 1000 runs was recorded. Fewer threads per block results in faster execution times on the GPUs.

|  | 128x128 | 256x256 | 512x512 |
|---|---|---|---|
| Python, CPU | 274 | 177 | 149 |
| Julia, CPU-serial (unoptimized) | 20 | 12 | 12 |
| Julia, CPU-serial (optimized) | 1 | 1 | 1 |
| Julia, GPU | **229** | **366** | **386** |

**Table 2.** Speed-up (bold) or slow-down (non-bold) factor compared to the optimized CPU-serial Julia version at the same resolution. GPU speed-ups are based on transferring arrays between GPU and CPU every ten time steps.

block size does have an effect on how quickly the kernel runs, so we benchmarked the evaluation time of the same kernel run with different block sizes, as shown in Figure 3. Smaller block sizes run faster on the GPUs by 15%. This is interesting to note, but GPU compute time is so small compared to the memory transfer time that thread tuning has little impact on the overall simulation time.

### 3.3 Julia-MPI versus Fortran-MPI

Julia and Fortran codes were compared on multi-node CPU clusters, where both used MPI for communication between processors. Comparisons were made with domains of 128, 256, and 512-squared grid cells solving the shallow water equations. All timing tests were conducted for 10 time steps and repeated 12 times on each processor count, spanning 2 to 2048 processors by powers of two. The vertical dimension included 100 layers to mimic ocean model arrays and provide sufficient computational work on each processor. Separate timers report on computational work versus MPI communication within the time-stepping routine. The i/o, initialization, and finalization time is excluded.

Simulations were conducted on Cori-Haswell at the National Energy Research Scientific Computing Center (NERSC). Cori-Haswell consists of 2,388 nodes in 14 cabinets, using Intel Xeon Processor E5-2698 v3 with a clock rate of 2.3 GHz. Each processor has 32 physical threads per node and two hyper-threads per core, with 128 GB of memory per node. The interconnect is a Cray Aries with Dragonfly topology and $> 45$ TB/s global peak bisection bandwidth. The Julia-MPI and Fortran-MPI tests were both run with up to 32 ranks per node.

The scaling plots in Figure 4 show that the Julia-MPI and Fortran-MPI models have identical performance at two cores; Julia-MPI is faster by up to a factor of two for mid-range core counts; and Fortran-MPI is 2x faster than Julia-MPI at higher ranges, depending on the resolution. For both languages, computation scales well with processor count, while communication does not, and communication progressively requires a much larger fraction of time at higher processor counts (Figure 5). Once computations are optimized, communication, which is fixed by the interconnect speed, will remain a bottleneck regardless of the language (see, e.g. Koldunov et al. (2019)). At the lowest resolution of 128x128, there is insufficient work beginning at 512 processors (which corresponds to 32 grid-cells per processor), and timing is dominated by communication, resulting in poor scaling above 512 processors. Communication times in Julia are much more variable than in Fortran across samples and

330 processor counts, as shown in the right column of Figure 4. When measuring computation time without communication (Figure 4, right column), Julia-MPI scales nearly perfectly, while Fortran-MPI computational time drops off from perfect scaling at 8 and 16 cores. This produces the Julia times that are 2x faster for the total times for mid-range processor counts of 16 and higher. Overall, Julia performance on CPU clusters is competitive with Fortran. Once the high-level codes have been optimized, the "winner" between Julia and Fortran will likely depend on the details of the MPI libraries and hardware.

## 4 Optimization Tips for Julia Developers

Julia serves the dual purpose of a prototyping language as well as a production language. Not only can we construct quick-to-write but slow-performing code (although still significantly faster than other development languages, as we saw with comparison to python) to demonstrate an idea, we can also spend a bit more time to carefully construct an optimized code to achieve performance on par with Fortran. Julia's ability to act as a prototyping language can be attributed to one of its key features: 340 dynamic typing. Just like Python, variables may be initialized without defining their types. However, Julia is also endowed with a static typing feature, even though it is optional. If the variable types are statically defined in a concrete fashion, performance is greatly improved. Julia activates its dynamic typing feature with an "Any" type which could be any type at run time. So, Julia must compile parts of the code on the fly (Julia Development Team, b). A method involving an "Any" type is compiled at run time for whatever type is actually provided during execution (called just-in-time compiling). The implication is that 345 without static typing, performance will greatly suffer from compilation at run time. Additionally, with concrete types, the Julia compiler may optimize the code much further than if it is compiled for an unknown type.

When first creating the MPAS shallow water core in Julia, we did not specify the array types, and let Julia assign them the "Any" type:

```
struct MPAS_Ocean
    layerThickness
    normalVelocity
    ...
end
```

However, by concretely defining these variables to be floating point arrays, we gain a substantial performance boost:

```
struct MPAS_Ocean
    layerThickness::Array{Float64}
    normalVelocity::Array{Float64}
end
```

When parallelizing for the graphics card, a different array type is used that is suited for GPUs. We tried defining an abstract 360 array type that encompasses both the CPU and GPU data types, so that CUDA.CuArrays and regular Arrays could be used
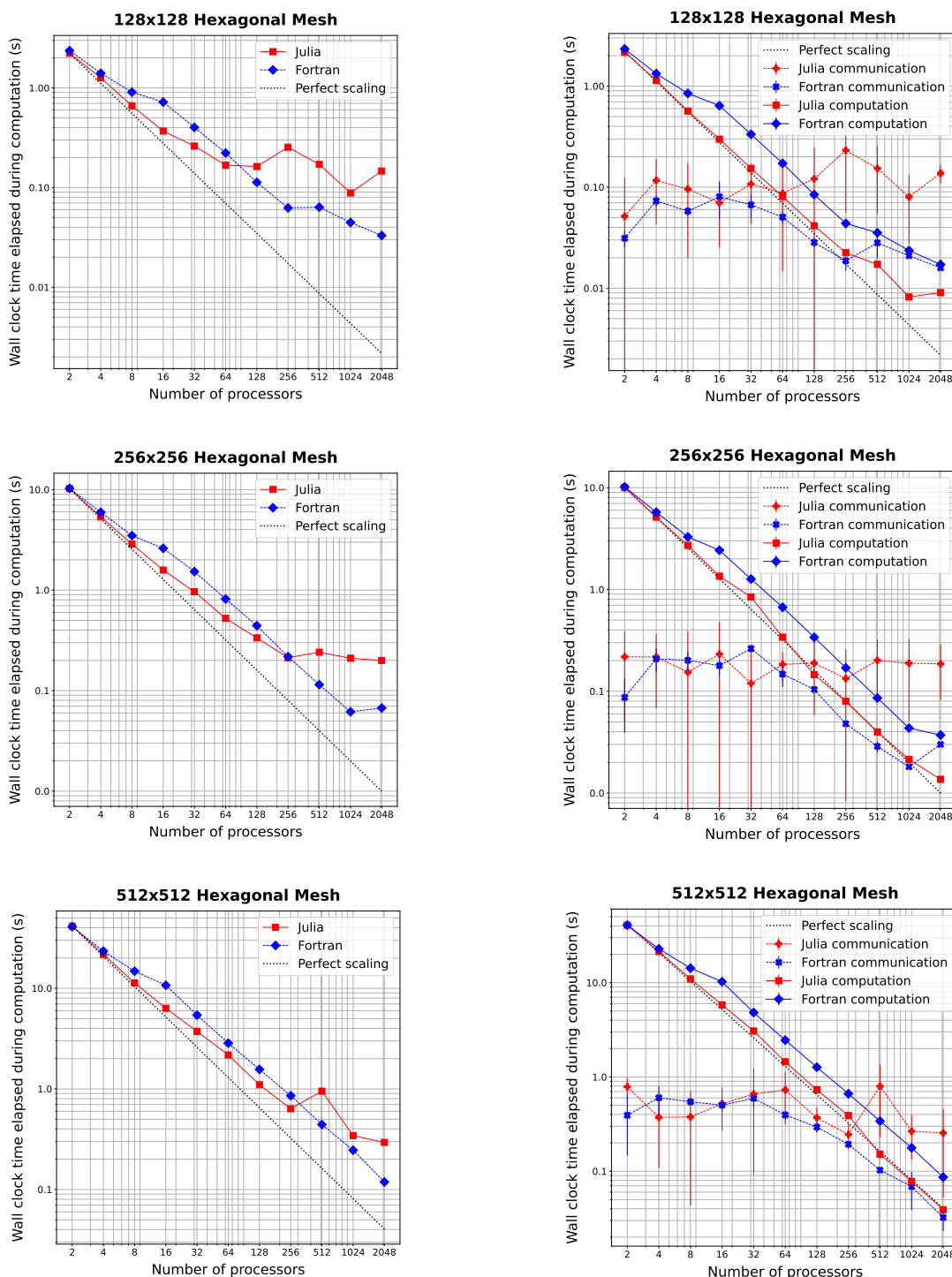
14

**Figure 4.** Wall clock time versus the number of processors to simulate 10 steps of the coastal Kelvin wave test with 100 layers. Left column shows total time without i/o; right column splits MPI communication and computation. Vertical lines display the standard deviation of communication times.
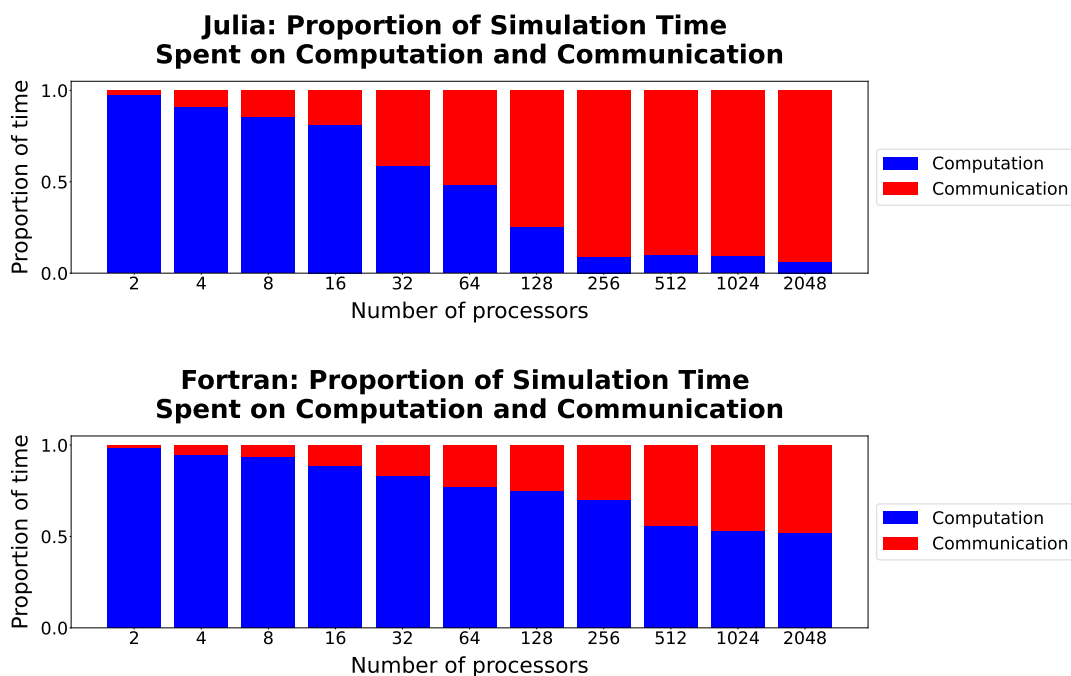
**Figure 5.** Comparison of the proportion of time spent in computation (blue) versus communication (red) in Julia-MPI (top) and Fortran-MPI (bottom) on the 128x128 hexagonal mesh. The relative time spent in communication increases dramatically at high processor counts.

interchangeably, allowing the model to be run on the GPU or CPU at will. We also used an abstract type specification on the contents of these arrays `F <: Float`, meaning any type extending the abstract floating point type can be used at runtime.

```
struct MPAS_Ocean{F <: AbstractFloat}
    layerThickness::AbstractArray{F}
    normalVelocity::AbstractArray{F}
end
```

365

This approach seems like it should be performant, since the types are defined before run time. However, abstract types, like an Any type, slow down execution since at run time they may actually be a different type that extends the abstract type (`CUDA.CuArray` or `Array`), meaning the compiler is doing just-in-time compiling. Similarly, specifying an inexact element type (`F <: AbstractFloat`) rather than a concrete type (`Float64`) is very inefficient.

370

Instead, two separate structures should be defined concretely when running on GPUs versus CPUs:

```
struct MPAS_Ocean_CUDA
    layerThickness::CUDA.CuArray{Float64,2}
    normalVelocity::CUDA.CuArray{Float64,2}
end
```

375

16

```
struct MPAS_Ocean
    layerThickness::Array{Float64,2}
    normalVelocity::Array{Float64,2}
end
```

380

Now the array types are concrete, element types are concrete (`Float64`), and the number of dimensions is specified (`Float64,2`). This code no longer has the advantageous feature of being able to switch between running on the CPU and GPU on the fly. However, the execution speed is massively improved. We found that making this change from abstract to concrete array types sped up computation by a factor of 34x.

385     The key in optimizing Julia code, we found, was reducing allocations. Memory allocation significantly slows down execution. And it is not always obvious what seemingly innocent actions may allocate memory. For example, simply reading a pair of values from an array with two columns:

```
cell1Index, cell2Index = cellsOnEdge[:,iEdge]
```

can allocate significant memory. In one test, this one line (executed repeatedly throughout the simulation) allocated 408 KiB.

390     This is because the line is really creating a tuple, not directly reading each column into the two scalar variables. If we separate this into two lines to enforce only using scalars and not allocating tuples or arrays:

```
cell1Index = cellsOnEdge[1,iEdge]
cell2Index = cellsOnEdge[2,iEdge]
```

then this cuts allocations to zero—making this line almost instantaneous, and dropping the time spent on the whole tendency

395     calculation from 198 $\mu s$ to 99 $\mu s$. That means this line alone was responsible for about 50% of the computation time, when it could be rewritten to take no time at all.

There are likely many inconspicuous lines like this lurking in one's Julia code, slowing it down substantially. Additionally, even one overlooked field which is not concretely typed may significantly slow execution. Luckily, Julia is equipped with a tool to quickly locate such memory-hoarding lines. This tool is called `@code_warntype`. Prefixing a function call with it

400     will print out a color-coded list breaking each line down to individual memory operations:

```
@code_warntype calculate_normal_velocity_tendency!(mpas)
```

It helpfully highlights inexact types and memory allocations with red, pointing a user right to the lines and fields that need to be optimized. This feature alone makes Julia very powerful for high-performance applications, significantly speeding up development time to optimize a model's performance.

405     Another very helpful tool when optimizing Julia code is `--track-allocations`, a command line option that can be added to any Julia execution as follows:

```
$ julia --track-allocations=user ./anyJuliascript.jl
```

A new file is created at `./anyJuliascript.jl.XXX.mem` (where XXX is some unique number). This file contains each line of the script prefixed by the number of memory allocations created by that line, giving a line-by-line breakdown of where

410    allocations occur.

## 5   Conclusions

As new programming languages and libraries become available, it is important for model developers to learn new techniques and evaluate them against their current methods. This is particularly true as computing architectures continue to evolve, and long-standing languages such as C++ and Fortran require additional libraries to remain competitive on new supercomputers.

415    In this work, we created three implementations of a shallow water model in Julia in order to compare ease of development and performance to standard Fortran and Python implementations. The three Julia codes were designed for single-CPU, GPU-enhanced single CPU, and parallelized multi-core CPU architectures. Julia-MPI speeds were identical to Fortran-MPI at low core counts, 2x faster for mid-range, and 2x slower at higher core counts. Julia-MPI exhibited better scaling than Fortran-MPI for computation-only times, and more variability for communication times.

420    The most surprising result of this study was the speed of computations on the GPUs—a speed-up of 40,000 to over 100,000 times compared to the CPU. Of course, this comes with the caveat that memory transfer between CPU and GPU can take thousands of times longer than the computation, up to 0.5s at our highest resolution. So the key is to transfer memory to and from the GPU as little as possible, which is a well-known practice. If one can fit the full resolution of a computational physics domain within the memory of a single graphics card and sample results rarely, GPUs offer significant speed-ups. For climate

425    models, a single low-resolution component may well fit into GPU memory if the developers are careful with their memory footprint. The difficulty is that including ocean, atmosphere, land, and sea ice components requires the use of multiple nodes, and inter-node communication will keep the model slow, regardless of the GPU speed. Higher-resolution domains will need many nodes for each component and present the same problem.

The shallow water equations are simple enough for rapid development and verification, yet contain the salient features of any

430    ocean model: intensive computation of the tendency terms, a time-stepping routine, and for the parallel version, interleaved halo communication of the partition boundary. Indeed, this layout, and the lessons learned here, apply to almost all computational physics codes.

This work specifically tests unstructured horizontal meshes, as opposed to structured quadrilateral grids. Unstructured meshes refer to a neighbor's index using additional pointer arrays, so require an extra memory access for horizontal stencils.

435   In structured grids, the physical neighbors are also neighbors in array space ($i + 1, j + 1$, etc), which leads to more contiguous memory access patterns that are easier for compilers to optimize. Our results show that unstructured meshes do not present any significant challenge in either Fortran or Julia. The use of a structured vertical index in the innermost position and testing with 100 layers provides sufficient contiguous memory access for cache locality.

In the end, we were impressed by our experience with Julia. It did fulfill the promise of fast and convenient prototyping, with

440    the ability to eventually run at high speeds on multiple high performance architectures—after some effort and lessons learned

by the developers. The Julia libraries for MPI and CUDA were powerful and convenient. E3SM does not have plans to develop model components with Julia, but this study provides a useful comparison to our C++ and Fortran codes as we move towards heterogeneous, exascale computers.

*Code and data availability.* Three code repositories were used for the performance comparisons in this study. These are publicly available on both GitHub and Zenodo:

1. Julia Shallow Water code for serial CPU, CUDA-GPU, and MPI-parallelized CPU

    GitHub: https://github.com/robertstrauss/MPAS_Ocean_Julia (license: GNU General Public License v3.0)

    Zenodo: https://doi.org/10.5281/zenodo.7493065 (license: Creative Commons Attribution 4.0 International)

2. Python Rotating Shallow Water Verification Suite

    GitHub: https://github.com/siddharthabishnu/Rotating_Shallow_Water_Verification_Suite.git. (license: LANL/UCAR*)
    This study used the specific code version https://github.com/siddharthabishnu/Rotating_Shallow_Water_Verification_Suite/tree/v1.0.1 (license: LANL/UCAR*)

    Zenodo: https://doi.org/10.5281/zenodo.7425628 (license: BSD 3-Clause "New" or "Revised")

3. Fortran-MPI MPAS Shallow Water code with Coastal Kelvin wave initial condition (Petersen et al., 2022)

    GitHub: https://github.com/MPAS-Dev/MPAS-Model. (license: LANL/UCAR*) This study used the specific code version https://github.com/mark-petersen/MPAS-Model/releases/tag/SW_julia_comparison_V1.0.

    Zenodo: https://doi.org/10.5281/zenodo.7439134 (license: Creative Commons Attribution 4.0 International)

The planar hexagonal MPAS-Ocean meshes used in this study for the numerical simulations and convergence tests of the coastal Kelvin wave and the inertia-gravity wave can be obtained from the Zenodo release of the Python Rotating Shallow Water Verification Suite Meshes at https://doi.org/10.5281/zenodo.7419817.

\* Code bases use the license found at https://github.com/MPAS-Dev/MPAS-Model/blob/master/LICENSE.

*Author contributions.* Code development, testing, and timing were conducted by RRS for Julia, SB for python, and MRP for Fortran. RRS led data analysis, plot generation, and Julia optimization. SB led the test case design and verification. The manuscript was written cooperatively by all authors. MRP conceptualized the project and conducted funding acquisition.

*Competing interests.* The authors declare no competing interests

# References

Bishnu, S.: Time-Stepping Methods for Partial Differential Equations and Ocean Models, Ph.D. thesis, Florida State University, https://doi.org/10.5281/zenodo.7439539, 2021.

Bishnu, S.: Rotating Shallow Water Verification Suite, https://doi.org/10.5281/zenodo.7425628, 2022.

480 Bishnu, S., Petersen, M., Quaife, B., and Schoonover, J.: Verification Suite of Test Cases for the Barotropic Solver of Ocean Models, https://doi.org/10.22541/essoar.167100170.03833124/v1, 2022.

Bleichrodt, F., Bisseling, R. H., and Dijkstra, H. A.: Accelerating a barotropic ocean model using a GPU, OCEAN MODEL, 41, 16–21, https://doi.org/10.1016/j.ocemod.2011.10.001, 2012.

Caldwell, P. M., Mametjanov, A., Tang, Q., Van Roekel, L. P., Golaz, J. C., et al.: The DOE E3SM Coupled Model Version 1: Description
485 and Results at High Resolution, J ADV MODEL EARTH SY, 11, 4095–4146, https://doi.org/10.1029/2019MS001870, 2019.

Cushman-Roisin, B. and Beckers, J.-M.: Introduction to geophysical fluid dynamics: physical and numerical aspects, Academic press, 2011.

Dalcín, L., Paz, R., and Storti, M.: MPI for Python, J PARALLEL DISTR COM, 65, 1108–1115, 2005.

Dalcín, L., Paz, R., Storti, M., and D'Elía, J.: MPI for Python: Performance improvements and MPI-2 extensions, J PARALLEL DISTR COM, 68, 655–662, 2008.

490 Gevorkyan, M. N., Demidova, A. V., Korolkova, A. V., and Kulyabov, D. S.: Statistically significant performance testing of Julia scientific programming language, J PHYS CONF SER, 1205, 012 017, https://doi.org/10.1088/1742-6596/1205/1/012017, 2019.

Golaz, J.-C., Caldwell, P. M., Van Roekel, L. P., Petersen, M. R., et al.: The DOE E3SM Coupled Model Version 1: Overview and Evaluation at Standard Resolution, J ADV MODEL EARTH SY, 11, 2089–2129, https://doi.org/10.1029/2018MS001603, 2019.

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
495 to Gpus Using OpenAcc, IEEE ACCESS, 7, 154 490–154 501, https://doi.org/10.1109/ACCESS.2019.2932443, 2019.

Julia Development Team: Introduction to CUDA, https://cuda.juliagpu.org/stable/tutorials/introduction/#A-simple-example-on-the-CPU, 2022, visited on 2022-12-13, a.

Julia Development Team: Eval of Julia code, https://docs.julialang.org/en/v1/devdocs/eval/#, 2016, visited on 2022-12-13, b.

Klöckner, A., Pinto, N., Lee, Y., Catanzaro, B., Ivanov, P., and Fasih, A.: PyCUDA and PyOpenCL: A Scripting-Based Approach to GPU
500 Run-Time Code Generation, PARALLEL COMPUT, 38, 157–174, https://doi.org/10.1016/j.parco.2011.09.001, 2012.

Klöwer, M., Hatfield, S., Croci, M., Düben, P. D., and Palmer, T. N.: Fluid Simulations Accelerated With 16 Bits: Approaching 4x Speedup on A64FX by Squeezing ShallowWaters.jl Into Float16, J ADV MODEL EARTH SY, 14, https://doi.org/10.1029/2021MS002684, 2022.

Koldunov, N. V., Aizinger, V., Rakowsky, N., Scholz, P., Sidorenko, D., Danilov, S., and Jung, T.: Scalability and some optimization of the Finite-volumE Sea ice–Ocean Model, Version 2.0 (FESOM2), GEOSCI MODEL DEV, 12, 3991–4012,
505 https://doi.org/10.5194/gmd-12-3991-2019, 2019.

Lam, S. K., Pitrou, A., and Seibert, S.: Numba: A llvm-based python jit compiler, in: Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC, pp. 1–6, 2015.

Lin, W.-C. and McIntosh-Smith, S.: Comparing Julia to Performance Portable Parallel Programming Models for HPC, in: 2021 International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS), pp. 94–105, IEEE,
510 St. Louis, MO, USA, https://doi.org/10.1109/PMBS54543.2021.00016, 2021.

Mielikainen, J., Huang, B., Huang, H.-L. A., and Goldberg, M. D.: Improved GPU/CUDA Based Parallel Weather and Research Forecast (WRF) Single Moment 5-Class (WSM5) Cloud Microphysics, IEEE J SEL TOP APPL, 5, 1256–1265, https://doi.org/10.1109/JSTARS.2012.2188780, 2012.

Perkel, J. M.: Julia: come for the syntax, stay for the speed, NATURE, 572, 141–142, https://doi.org/10.1038/d41586-019-02310-3, 2019.

515  Petersen, M. R., Jacobsen, D. W., Ringler, T. D., Hecht, M. W., and Maltrud, M. E.: Evaluation of the Arbitrary Lagrangian–Eulerian Vertical Coordinate Method in the MPAS-Ocean Model, OCEAN MODEL, 86, 93–113, https://doi.org/10.1016/j.ocemod.2014.12.004, 2015.

Petersen, M. R., Asay-Davis, X. S., Berres, A. S., Chen, Q., Feige, N., Hoffman, M. J., Jacobsen, D. W., Jones, P. W., Maltrud, M. E., Price, S. F., Ringler, T. D., Streletz, G. J., Turner, A. K., Van Roekel, L. P., Veneziani, M., Wolfe, J. D., Wolfram, P. J., and Woodring, J. L.: An Evaluation of the Ocean and Sea Ice Climate of E3SM Using MPAS and Interannual CORE-II Forcing, J ADV MODEL EARTH SY, 11,
520  1438–1458, https://doi.org/10.1029/2018MS001373, 2019.

Petersen, M. R., Bishnu, S., and Strauss, R. R.: MPAS-Ocean Shallow Water Performance Test Case, https://doi.org/10.5281/zenodo.7439134, 2022.

Ramadhan, A., Wagner, G. L., Hill, C., Campin, J.-M., Churavy, V., Besard, T., Souza, A., Edelman, A., Ferrari, R., and Marshall, J.: Oceananigans.jl: Fast and friendly geophysical fluid dynamics on GPUs, J. OPEN SOURCE SOFTW., 5, 2018,
525  https://doi.org/10.21105/joss.02018, 2020.

Ringler, T. D., Thuburn, J., Klemp, J. B., and Skamarock, W. C.: A unified approach to energy conservation and potential vorticity dynamics for arbitrarily-structured C-grids, J COMPUT PHYS, 229, 3065–3090, 2010.

Ringler, T. D., Petersen, M. R., Higdon, R. L., Jacobsen, D., Jones, P. W., and Maltrud, M.: A multi-resolution approach to global ocean modeling, OCEAN MODEL, 69, 211–232, 2013.

530  Shchepetkin, A. F. and McWilliams, J. C.: The regional oceanic modeling system (ROMS): a split-explicit, free-surface, topography-following-coordinate oceanic model, OCEAN MODEL, 9, 347–404, 2005.

Srinath, A.: Accelerating Python on GPUs with nvc++ and Cython, https://developer.nvidia.com/blog/accelerating-python-on-gpus-with-nvc-and-cython/, 2020, visited on 2022-12-13.

Strauss, R. R.: Julia Layered Shallow Water Model on Various Hardwares, https://doi.org/10.5281/zenodo.7493065, 2023.

535  Thuburn, J., Ringler, T. D., Skamarock, W. C., and Klemp, J. B.: Numerical representation of geostrophic modes on arbitrarily structured C-grids, J COMPUT PHYS, 228, 8321–8335, 2009.

Trott, C. R., Lebrun-Grandié, D., et al.: Kokkos 3: Programming Model Extensions for the Exascale Era, IEEE T PARALL DISTR, 33, 805–817, https://doi.org/10.1109/TPDS.2021.3097283, 2022.

Xu, S., Huang, X., Zhang, Y., Hu, Y., and Yang, G.: A customized GPU acceleration of the princeton ocean model,
540  in: 2014 IEEE 25th International Conference on Application-Specific Systems, Architectures and Processors, pp. 192–193, https://doi.org/10.1109/ASAP.2014.6868661, 2014.

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, GEOSCI MODEL DEV, 8, 2815–2827, https://doi.org/10.5194/gmd-8-2815-2015, 2015.

Ye, Y., Song, Z., Zhou, S., Liu, Y., Shu, Q., Wang, B., Liu, W., Qiao, F., and Wang, L.: swNEMO_v4.0: an ocean model based on NEMO4 for
545  the new-generation Sunway supercomputer, GEOSCI MODEL DEV, 15, 5739–5756, https://doi.org/10.5194/gmd-15-5739-2022, 2022.

Zhao, X.-d., Liang, S.-x., Sun, Z.-c., Zhao, X.-z., Sun, J.-w., and Liu, Z.-b.: A GPU accelerated finite volume coastal ocean model, J HYDRODYN, Ser. B, 29, 679–690, https://doi.org/10.1016/S1001-6058(16)60780-1, 2017.