



# Unlocking the power of parallel computing: GPU technologies for ocean forecasting

Andrew R. Porter<sup>1</sup> and Patrick Heimbach<sup>2</sup>

<sup>1</sup>Science and Technology Facilities Council, Daresbury Laboratory, Hartree Centre, Daresbury, UK

<sup>2</sup>Oden Institute for Computational Engineering and Sciences, The University of Texas at Austin,  
Austin, Texas, USA

**Correspondence:** Andrew R. Porter ([andrew.porter@stfc.ac.uk](mailto:andrew.porter@stfc.ac.uk))

Received: 17 September 2024 – Discussion started: 26 September 2024

Revised: 4 February 2025 – Accepted: 10 February 2025 – Published: 2 June 2025

**Abstract.** Operational ocean forecasting systems (OOFs) are complex engines that must execute ocean models with high performance to provide timely products and datasets. Significant computational resources are then needed to run high-fidelity models, and, historically, the technological evolution of microprocessors has constrained data-parallel scientific computation. Today, graphics processing units (GPUs) offer a rapidly growing and valuable source of computing power rivaling the traditional CPU-based machines: the exploitation of thousands of threads can significantly accelerate the execution of many models, ranging from traditional HPC workloads of finite difference, finite volume, and finite element modelling through to the training of deep neural networks used in machine learning (ML) and artificial intelligence. Despite the advantages, GPU usage in ocean forecasting is still limited due to the legacy of CPU-based model implementations and the intrinsic complexity of porting core models to GPU architectures. This review explores the potential use of GPU in ocean forecasting and how the computational characteristics of ocean models can influence the suitability of GPU architectures for the execution of the overall value chain: it discusses the current approaches to code (and performance) portability, from CPU to GPU, including tools that perform code transformation, easing the adaptation of Fortran code for GPU execution (like PSyclone), the direct use of OpenACC directives (like ICON-O), the adoption of specific frameworks that facilitate the management of parallel execution across different architectures, and the use of new programming languages and paradigms.

## 1 Introduction

Operational ocean forecasting systems (OOFs) are computationally demanding, and large computing resources are required in order to run models of useful fidelity. However, this is a time of great upheaval in the development of computer architectures. The ever-shrinking size of transistors means that current leakage (and the resulting heat generated) now presents a significant challenge to chip designers. This breakdown of Dennard scaling (transistor power consumption is proportional to area as in Dennard et al., 1974) began in about 2006 and means that it is no longer straightforward to continually increase the clock frequency of processors. Historically, this has been the main source of performance improve-

ment from one generation of processor to the next (Fig. 1). Although the number of transistors per device continues to rise, they are increasingly being used to implement larger numbers of execution cores. It is then the job of the application to make use of these additional cores to achieve a performance improvement. Graphics processing units (GPUs) are a natural consequence of this evolution. Originally developed to accelerate the rendering of computer-generated images (a naturally data-parallel task thanks to the division of an image into pixels), scientists were quick to seize on their potential to accelerate data-parallel scientific computation. Therefore, manufacturers today produce HPC-specific “GPUs” that are purely intended for computation. The suitability of this hardware for the training of deep neural networks used in ma-

chine learning (ML) and artificial intelligence has stimulated massive development and competition amongst GPU vendors. Because of the exploding interest of AI applications in virtually all sectors of industry, the commercial HPC market is undergoing a seismic shift toward GPU-based hardware, with serious implications for available HPC architectures in the future, to which OOFs will have to adapt.

Unlike CPUs, which tend to have relatively few but powerful (general purpose) processor cores, GPUs support hundreds of simpler cores running thousands of threads which can obtain data from memory very efficiently. The simplicity of these cores makes them more energy-efficient; therefore GPUs tend to offer significantly greater performance per watt. With the energy consumption of large computing facilities now the key design criterion, GPUs are an important part of the technology being used in the push towards exascale performance and beyond (e.g. Draeger and Siegel, 2023). As an illustration, in the November 2024 incarnation of the TOP500 list (Strohmaier et al., 2024), 9 of the machines in the top 10 are equipped with GPU accelerators from NVIDIA, Intel, or AMD. Although CPUs are present in these machines, their primary role is to host the GPUs which provide the bulk of the computing performance. GPUs are therefore a major feature of the current HPC landscape, and their importance and pervasiveness are only set to increase.

## 2 Computational characteristics of ocean models

To understand why GPUs are well suited to running OOFs, it is important to consider their computational characteristics. The equations describing ocean evolution form a system of partial differential equations that are solved numerically by discretising the model domain and then using a finite difference, finite volume, or finite element scheme. In these forms, the bulk of the computational work takes the form of stencil computations, where the update of a field at a given grid location requires that many other field values be read from neighbouring locations. This means that the limiting factor in the rate at which these computations can be done is how quickly all these values can be fetched from memory (so-called “memory bandwidth”). Finite element schemes do have the advantage of shifting the balance in favour of doing more arithmetic operations, but memory bandwidth still tends to dominate. These computations are, of course, repeated across the entire model grid, meaning that it is a single instruction, multiple data (SIMD) problem. OOFs are therefore a very good fit for GPU architectures, which naturally support massively data-parallel problems and typically provide much higher memory bandwidth than CPUs.

For execution on distributed-memory computers, OOFs typically use a geographical domain decomposition where each processor is assigned a part of the model domain. In order to handle stencil updates at the boundaries of a processor’s sub-domain, it must exchange information with those

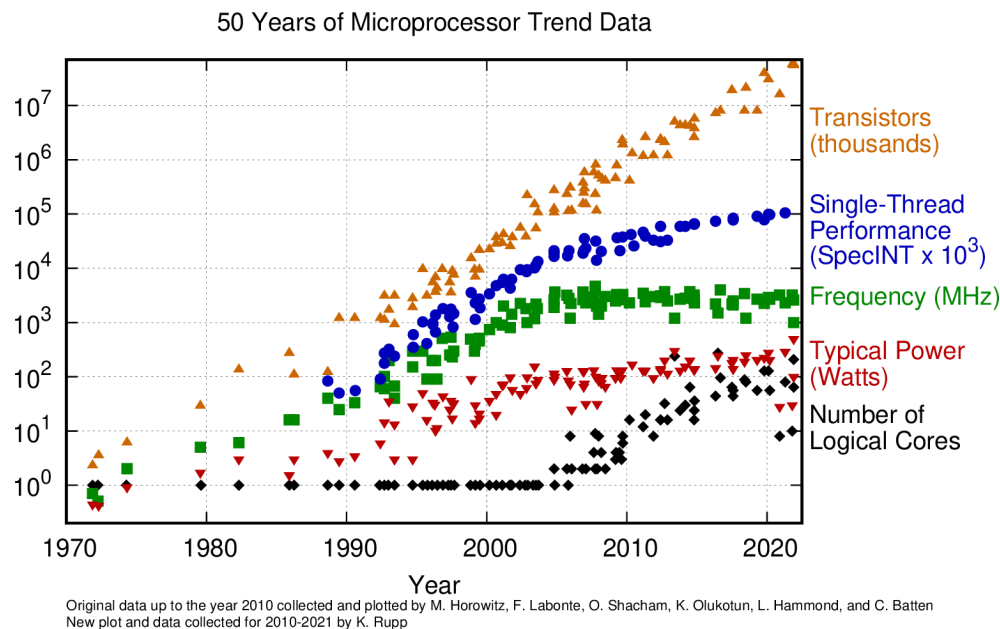
processors operating on neighbouring sub-domains. Obviously, there is a cost associated with performing these exchanges, which high-performance processor interconnects can only do so much to mitigate. As more processors are thrown at a problem in order to reduce the time to solution, the size of their sub-domains decreases and so does the amount of computation that each must perform. Consequently, the relative cost of inter-processor communication becomes more significant and, after a certain point (the “strong-scaling limit”), will begin to dominate. At this point, using further processors will bring only limited performance improvements, if any.

Inter-processor communication on a GPU-based machine can be more costly, as messages may have to go via the CPUs hosting the GPUs, unless a machine has both hardware and software support for direct GPU–GPU communication. Communication-avoiding/minimising strategies are therefore more important on these architectures. These can include algorithmic design (e.g. Silvestri et al., 2024) to allow the overlap of communication and computation or simply the use of wider halo regions to reduce the frequency of halo exchanges.

## 3 The use of GPUs in ocean forecasting

Although GPUs are now a well-established HPC technology with potentially significant performance advantages for OOFs, they are not yet widely adopted in the ocean-forecasting community. For example, in Europe, NEMO (Madec et al., 2024) is the most important ocean-modelling framework; it is used operationally by Mercator Ocean International, the European Centre for Medium-Range Weather Forecasts (ECMWF), the UK Met Office, the Euro-Mediterranean Center on Climate Change, and other institutes worldwide. NEMO is implemented in Fortran and parallelised with MPI and, as such, is limited to running on CPUs only. The German Weather Service (DWD) uses ICON-O (Korn, 2017), which is also a Fortran model. Experiments are in progress with the use of OpenACC directives to extend this code to make use of GPUs, but this functionality is not used operationally.

In the US, NOAA’s Real-Time Ocean Forecast System (<https://polar.ncep.noaa.gov/global/>, last access: 14 April 2025) is based on the Hybrid Coordinate Ocean Model (HYCOM; Chassignet et al., 2009). HYCOM is also a Fortran code, parallelised using a combination of OpenMP and MPI. Although not used operationally, the Energy Exascale Earth System Model (E3SM) is also significant. It utilises the ocean, sea ice, and land ice versions of the Model for Prediction Across Scales (MPAS; Ringler et al., 2013), which again is implemented in Fortran with MPI. Although a port of this was attempted through the addition of OpenACC directives, it has been abandoned due to poor GPU performance (Mark R. Petersen, personal communication, 2024). Instead,



**Figure 1.** The breakdown of Dennard scaling, shown by 50 years of microprocessor (CPU) evolution (Rupp, 2022).

a new ocean model on unstructured meshes named Omega is being developed in C++ from the ground up. Other widely used ocean general circulation models include the MIT General Circulation Model (MITgcm; Marshall et al., 1997) and the Modular Ocean Model, version 6 (MOM6; Adcroft et al., 2019), both of which again are Fortran codes with support for distributed- and shared-memory parallelism on CPU.

The Japanese Meteorological Agency runs operational forecasts using the Meteorological Research Institute Community Ocean Model (MRI.COM) (Tsujino et al., 2010). As with the previous models, this is also implemented in Fortran with MPI and thus only runs on CPU.

For regional (as opposed to global) forecasts, the Rutgers Regional Ocean Modeling System (ROMS) (Shchepetkin and McWilliams, 2003) is used by centres worldwide, including the Japan Fisheries Research and Education Agency, the Australian Bureau of Meteorology, and the Irish Marine Institute. ROMS is also a Fortran code parallelised using either MPI or OpenMP (but not both combined) and is thus restricted to CPU execution. Although various projects have ported the code to different architectures (including the Sunway architecture for China's Tianhe machine; Liu et al., 2019), these are all standalone pieces of work that have not made it back into the main code base.

## 4 Discussion

From the preceding section, it is clear that OOFs are currently largely implemented in Fortran with no or limited support for execution on GPU devices. The problem here is that OOFs comprise large and complex codes which typically

have a lifetime of decades and are constantly being updated with new science by multiple developers. Maintainability, allowing for the fact that the majority of developers will be specialists in their scientific domain rather than in HPC, is therefore of vital importance. Given that such codes are often shared between organisations, they must also run with good performance on different types of architecture (i.e. be “performance-portable”).

Previously, one generation of supercomputers looked much like the last; therefore the evolution of these computer models was not a significant problem. However, the proliferation of computer hardware (and, crucially, the programming models needed to target them) that has resulted from the breakdown of Dennard scaling has changed this (Balaji, 2021). With the average supercomputer having a lifetime of just some 5 years, OOFs are now facing the problem of adapting to future supercomputer architectures, and this is difficult because the aims of performance, performance portability, and code maintainability often conflict with each other (Lawrence et al., 2018).

*Transformation of existing codes.* To date there have been various approaches to this problem. NEMO v.5.0 (Madec et al., 2024) has adopted the PSyclone code transformation tool (Adams et al., 2019), which enables an HPC expert to transform Fortran source code such that it may be executed on GPUs using whichever programming model (i.e. OpenACC or OpenMP) is required. Previous, unpublished work found that, for a low-resolution ( $1^\circ$ ) global mesh, a single NVIDIA V100 GPU performed some  $3.6 \times$  better than an HPC-class Intel socket. For a high-resolution ( $1/12\text{th}^\circ$ ) global mesh,  $\sim 90$  A100 GPUs gave the same performance

as  $\sim 270$  Intel sockets. In both cases, this is an ocean-only configuration, with virtually all computing being performed on the GPUs. This is important, since any computation happening on the CPU incurs substantial data transfer costs as data are moved from the GPU to the CPU, updated, and then transferred back to the GPU. The advent of hardware support for unified CPU/GPU memory should reduce the cost of this. As noted earlier, ICON-O is being extended manually with OpenACC directives. There are examples of recent (i.e. experimental) models that have moved away from Fortran in favour of higher-level programming approaches. Thetis (Kärnä et al., 2018) implements a discontinuous Galerkin method for solving the 3D hydrostatic equations using the Firedrake framework. This permits the scientist to express their scheme in the Python implementation of Unified Form Language (Alnæs et al., 2014). The necessary code is then generated automatically. The Veros model (Häfner et al., 2021) takes a slightly different approach: its dynamical core is a direct Python translation of a Fortran code and thus retains explicit MPI parallelisation. The JAX system (<http://github.com/google/jax>, last access: 14 April 2025) for Python is then used to generate performant code for both CPU and GPU. The authors report that the Python version running on 16 A100 GPUs gives the same performance as 2000 CPU cores for the Fortran version (although this comparison is slightly unfair, as the CPUs used are several generations older than the GPUs).

**Performance portability tools.** Another popular approach to performance portability is to implement a model using a framework that takes care of parallel execution on a target platform. Frameworks such as Kokkos (Carter Edwards et al., 2014), SYCL, and OpenMP are good examples, and the new Omega ocean component of E3SM mentioned previously is being developed to use Kokkos. In principle, this approach retains single-source science code while enabling portability to a variety of different hardware. However, it is hard to insulate the oceanographer from the syntax of the framework (which is often only available in C++), and, while the framework may be portable, obtaining good performance often requires that it be used in a different way from one platform to another. In OpenMP, for instance, the directives needed to parallelise a code for a multi-core CPU are not the same as those needed to offload code to an accelerator.

**New programming languages.** The Climate Modeling Alliance (CliMA) has adopted a radically new approach by rewriting ocean and atmospheric models from scratch using the programming language Julia (Perkel, 2019; Sridhar et al., 2022). Designed to overcome the “two-language problem” (Churavy et al., 2022), Julia is ideally suited to harness emerging HPC architectures based on GPUs (Besard et al., 2017; Bezanson et al., 2017). First results with CliMA’s ocean model, Oceananigans.jl (Ramadhan et al., 2020), run on 64 NVIDIA A100 GPUs exhibit 10 simulated years per day (SYPD) at 8 km horizontal resolution (Silverstri et al., 2024). This performance is similar to current-

generation CPU-based ocean climate models run at much coarser resolution (order of 25–50 km resolution). Similarly promising benchmarks have been obtained with a barotropic configuration of a prototype of MPAS-Ocean, rewritten in Julia (Bishnu et al., 2023). Such performance gains hold great promise for accelerating operational ocean prediction at high spatial resolution run on emerging HPC hardware.

**Toward energy-efficient simulations.** Increased resolution, process representation, and data intensity in ocean and climate modelling is vastly expanding the need for compute cycles (more cores and smaller time steps). As a result, the ocean, atmosphere, and climate modelling community has recognised the need for their simulations to become more energy-efficient and to reduce their carbon footprint (Loft, 2020; Acosta et al., 2024; Voosen, 2024). Owing to their architecture, GPUs can play a significant role in reducing energy requirements. A related research frontier being spearheaded by the atmospheric modelling community is the use of mixed or reduced precision to speed up simulations (Freytag et al., 2022; Klöwer et al., 2022; Paxton et al., 2022), with a potentially desirable side effect of natively capturing stochastic parameterisations (Kimpson et al., 2023). GPUs are ideally suited for such approaches, but successful implementation depends heavily on the model’s numerical algorithms.

**Data-driven operational ocean forecasting.** Operational weather and ocean forecasting are facing the potential of a paradigm shift with the advent of powerful, purely data-driven methods. The numerical weather prediction (NWP) community has spearheaded the development of machine-learning-based emulators that perform several orders of magnitude faster than physics-based models (e.g. Bouallègue et al., 2024; Rasp et al., 2024). Such emulators have the potential to revolutionise probabilistic forecasting and uncertainty quantification, among others. The computational patterns underlying the ML algorithms, such as parallel matrix multiplication, are ideally suited for general-purpose GPU architectures. While these methods have been driven to a large extent by private sector entities and require access to increasingly large GPU-based HPC systems for training, corresponding efforts in operational ocean forecasting are only now beginning to catch up. A review of the rapidly changing landscape of AI methods in the context of ocean forecasting is attempted in Heimbach et al. (2025; in this report).

**Data availability.** No datasets were used in this article.

**Author contributions.** ARP created the first draft of this work. PH assisted with updating the text in light of the reviews received.



**Competing interests.** Andrew R. Porter has declared that he is an author of the PSyclone package and a co-chair of the NEMO HPC Working Group.

**Disclaimer.** Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims made in the text, published maps, institutional affiliations, or any other geographical representation in this paper. While Copernicus Publications makes every effort to include appropriate place names, the final responsibility lies with the authors.

**Acknowledgements.** The authors would like to thank the reviewers for bringing the E3SM work to their attention.

**Review statement.** This paper was edited by Jay Pearlman and reviewed by Mark R. Petersen and one anonymous referee.

## References

- Acosta, M. C., Palomas, S., Paronuzzi Ticco, S. V., Utrera, G., Biercamp, J., Bretonniere, P.-A., Budich, R., Castrillo, M., Caubel, A., Doblas-Reyes, F., Epicoco, I., Fladrich, U., Joussaume, S., Kumar Gupta, A., Lawrence, B., Le Sager, P., Lister, G., Moine, M.-P., Rioual, J.-C., Valcke, S., Zadeh, N., and Balaji, V.: The computational and energy cost of simulation and storage for climate science: lessons from CMIP6, *Geosci. Model Dev.*, 17, 3081–3098, <https://doi.org/10.5194/gmd-17-3081-2024>, 2024.
- Adams, S. V., Ford, R. W., Hambley, M., Hobson, J. M., Kavcic, I., Maynard, C. M., Melvin, T., Mueller, E. H., Mullerworth, S., Porter, A. R., Rezny, M., Shipway, B. J., and Wong, R.: LFRic: Meeting the challenges of scalability and performance portability in Weather and Climate models, *J. Parallel Distr. Com.*, 132, 383–396, <https://doi.org/10.1016/j.jpdc.2019.02.007>, 2019.
- Adcroft, A., Anderson, W., Balaji, V., Blanton, C., Bushuk, M., Dufour, C. O., Dunne, J. P., Griffies, S. M., Hallberg, R., Harrison, M. J., Held, I. M., Jansen, M. F., John, J. G., Krasting, J. P., Langenhorst, A. R., Legg, S., Liang, Z., McHugh, C., Radhakrishnan, A., Reichl, B. G., Rosati, T., Samuels, B. L., Shao, A., Stouffer, R., Winton, M., Wittenberg, A. T., Xiang, B., Zadeh, N., and Zhang, R.: The GFDL Global Ocean and Sea Ice Model OM4.0: Model Description and Simulation Features, *J. Adv. Model. Earth Sy.*, 11, 3167–3211, <https://doi.org/10.1029/2019ms001726>, 2019.
- Alnæs, M. S., Logg, A., Ølgaard, K. B., Rognes, M. E., and Wells, G. N.: Unified Form Language: A Domain-Specific Language for Weak Formulations of Partial Differential Equations, *ACM Trans. Math. Softw.*, 40, 1–37, <https://doi.org/10.1145/2566630>, 2014.
- Balaji, V.: Climbing down Charney's ladder: machine learning and the post-Dennard era of computational climate science, *Philos. T. Roy. Soc. A*, 379, 20200085, <https://doi.org/10.1098/rsta.2020.0085>, 2021.
- Besard, T., Foket, C., and Sutter, B. D.: Effective Extensible Programming: Unleashing Julia on GPUs, *IEEE T. Parall. Distr.*, 30, 827–841, <https://doi.org/10.1109/tpds.2018.2872064>, 2017.
- Bezanson, J., Edelman, A., Karpinski, S., and Shah, V. B.: Julia: A Fresh Approach to Numerical Computing, *SIAM Rev.*, 59, 65–98, <https://doi.org/10.1137/141000671>, 2017.
- Bishnu, S., Strauss, R. R., and Petersen, M. R.: Comparing the Performance of Julia on CPUs versus GPUs and Julia-MPI versus Fortran-MPI: a case study with MPAS-Ocean (Version 7.1), *Geosci. Model Dev.*, 16, 5539–5559, <https://doi.org/10.5194/gmd-16-5539-2023>, 2023.
- Bouallègue, Z. B., Clare, M. C. A., Magnusson, L., Gascón, E., Maier-Gerber, M., Janoušek, M., Rodwell, M., Pinault, F., Damsch, J. S., Lang, S. T. K., Raoult, B., Rabier, F., Chevallier, M., Sandu, I., Dueben, P., Chantry, M., and Pappenberger, F.: The Rise of Data-Driven Weather Forecasting: A First Statistical Assessment of Machine Learning-Based Weather Forecasts in an Operational-Like Context, *B. Am. Meteorol. Soc.*, 105, E864–E883, <https://doi.org/10.1175/bams-d-23-0162.1>, 2024.
- Carter Edwards, H., Trott, C. R., and Sunderland, D.: Kokkos: enabling manycore performance portability through polymorphic memory access patterns, *J. Parallel Distr. Com.*, 74, 3202–3216, <https://www.osti.gov/servlets/purl/1106586> (last access: 14 April 2025), 2014.
- Chassignet, E. P., Hurlburt, H. E., Metzger, E. J., Smedstad, O. M., Cummings, J., Halliwell, G. R., Bleck, R., Baraille, R., Wallcraft, A. J., Lozano, C., Tolman, H. L., Srinivasan, A., Hankin, S., Cornillon, P., Weisberg, R., Barth, A., He, R., Werner, F., and Wilkin, J.: U.S. GODAE: Global Ocean Prediction with the HYbrid Coordinate Ocean Model (HYCOM), *Oceanography*, 22, 64–75, <https://doi.org/10.5670/oceanog.2009.39>, 2009.
- Churavy, V., Godoy, W. F., Bauer, C., Ranocha, H., Schlottke-Lakemper, M., Räss, L., Blaschke, J., Giordano, M., Schnetter, E., Omlin, S., Vetter, J. S., and Edelman, A.: Bridging HPC Communities through the Julia Programming Language, *arXiv [preprint]*, <https://doi.org/10.48550/arxiv.2211.02740>, 4 November 2022.
- Dennard, R. H., Gaensslen, F., Yu, H., Rideout, L., Bassous, E., and LeBlanc, A.: Design of ion-implanted MOSFET's with very small physical dimensions, *IEEE J. Solid-St. Circ.*, SC-9, 256–268, 1974.
- Draeger, E. W. and Siegel, A.: Exascale Was Not Inevitable; Neither Is What Comes Next, *Comput. Sci. Eng.*, 25, 79–83, <https://doi.org/10.1109/mcse.2023.3311375>, 2023.
- Freytag, G., Lima, J. V. F., Rech, P., and Navaux, P. O. A.: Impact of Reduced and Mixed-Precision on the Efficiency of a Multi-GPU Platform on CFD Applications, *Lect. Notes Comput. Sc.*, 13380, 570–587, [https://doi.org/10.1007/978-3-031-10542-5\\_39](https://doi.org/10.1007/978-3-031-10542-5_39), 2022.
- Häfner, D., Nuterman, R., and Jochum, M.: Fast, cheap, and turbulent – Global ocean modeling with GPU acceleration in Python, *J. Adv. Model. Earth Sy.*, 13, e2021MS002717, <https://doi.org/10.1029/2021MS002717>, 2021.
- Heimbach, P., O'Donncha, F., Smith, T., Garcia-Valdecasas, J. M., Arnaud, A., and Wan, L.: Crafting the Future: Machine Learning for Ocean Forecasting, in: *Ocean prediction: present status and state of the art (OPSR)*, edited by: Álvarez Fanjul, E., Ciliberti, S. A., Pearlman, J., Wilmer-Becker, K., and Behera, S., Copernicus Publications, State Planet, 5-opsr, 22, <https://doi.org/10.5194/sp-5-opsr-22-2025>, 2025.
- Kärnä, T., Kramer, S. C., Mitchell, L., Ham, D. A., Piggott, M. D., and Baptista, A. M.: Thetis coastal ocean model: discontinuous Galerkin discretization for the three-dimensional

- hydrostatic equations, *Geosci. Model Dev.*, 11, 4359–4382, <https://doi.org/10.5194/gmd-11-4359-2018>, 2018.
- Kimpson, T., Paxton, E. A., Chantry, M., and Palmer, T.: Climate-change modelling at reduced floating-point precision with stochastic rounding, *Q. J. Roy. Meteor. Soc.*, 149, 843–855, <https://doi.org/10.1002/qj.4435>, 2023.
- 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, e2021MS002684, <https://doi.org/10.1029/2021ms002684>, 2022.
- Korn, P.: Formulation of an unstructured grid model for global ocean dynamics, *J. Comput. Phys.*, 339, 525–552, <https://doi.org/10.1016/j.jcp.2017.03.009>, 2017.
- Lawrence, B. N., Rezny, M., Budich, R., Bauer, P., Behrens, J., Carter, M., Deconinck, W., Ford, R., Maynard, C., Mullerworth, S., Osuna, C., Porter, A., Serradell, K., Valcke, S., Wedi, N., and Wilson, S.: Crossing the chasm: how to develop weather and climate models for next generation computers?, *Geosci. Model Dev.*, 11, 1799–1821, <https://doi.org/10.5194/gmd-11-1799-2018>, 2018.
- Liu, T., Zhuang, Y., Tian, M., Pan, J., Zeng, Y., Guo, Y., and Yang, M.: Parallel Implementation and Optimization of Regional Ocean Modeling System (ROMS) Based on Sunway SW26010 Many-Core Processor, *IEEE Access*, 7, 146170–146182, <https://doi.org/10.1109/ACCESS.2019.2944922>, 2019.
- Loft, R.: Earth System Modeling Must Become More Energy Efficient, *Eos, Transactions American Geophysical Union*, 101, <https://doi.org/10.1029/2020eo147051>, 2020.
- Madec, G., Bell, M., Benshila, R., Blaker, A., Boudrallé-Badie, R., Bricaud, C., Bruciaferri, D., Carneiro, D., Castrillo, M., Calvert, D., Chanut, J., Clementi, E., Coward, A., de, C., Dobricic, S., Epicoco, I., Éthé, C., Fiedler, E., Ford, D., Furner, R., Ganderton, J., Graham, T., Harle, J., Hutchinson, K., Iovino, D., King, R., Lea, D., Levy, C., Lovato, T., Maisonnave, E., Mak, J., Manuel, J., Martin, M., Martin, N., Martins, D., Masson, S., Mathiot, P., Mele, F., Mocavero, S., Moulin, A., Müller, S., Nurser, G., Oddo, P., Paronuzzi, S., Paul, J., Peltier, M., Person, R., Rousset, C., Rynders, S., Samson, G., Schroeder, D., Storkey, D., Storto, A., Téchené, S., Vancoppenolle, M., and Wilson, C.: NEMO Ocean Engine Reference Manual (v5.0), Zenodo [documentation], <https://doi.org/10.5281/zenodo.14515373>, 2024.
- Marshall, J., Adcroft, A., Hill, C., Perelman, L., and Heisey, C.: A finite-volume, incompressible Navier Stokes model for studies of the ocean on parallel computers, *J. Geophys. Res.*, 102, 5753–5766, <https://doi.org/10.1029/96JC02775>, 1997.
- Paxton, E. A., Chantry, M., Klöwer, M., Saffin, L., and Palmer, T.: Climate Modeling in Low Precision: Effects of Both Deterministic and Stochastic Rounding, *J. Climate*, 35, 1215–1229, <https://doi.org/10.1175/jcli-d-21-0343.1>, 2022.
- Perkel, J. M.: Julia: come for the syntax, stay for the speed, *Springer Nature*, 141–142, <http://www.nature.com/articles/d41586-019-02310-3> (last access: 14 April 2025), 2019.
- Ramadhan, A., Wagner, G., 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, *Journal of Open Source Software*, 5, 2018, <https://doi.org/10.21105/joss.02018>, 2020.
- Rasp, S., Hoyer, S., Merose, A., Langmore, I., Battaglia, P., Russell, T., Sanchez-Gonzalez, A., Yang, V., Carver, R., Agrawal, S., Chantry, M., Bouallegue, Z. B., Dueben, P., Bromberg, C., Sisk, J., Barrington, L., Bell, A., and Sha, F.: WeatherBench 2: A Benchmark for the Next Generation of Data-Driven Global Weather Models, *J. Adv. Model. Earth Sy.*, 16, e2023MS004019, <https://doi.org/10.1029/2023ms004019>, 2024.
- Ringler, T. Petersen, M., Higdon, R. L., Jacobsen, D., Jones, P. W., and Maltrud, M.: A multi-resolution approach to global ocean modeling, *Ocean Model.*, 69, 211–232, <https://doi.org/10.1016/j.ocemod.2013.04.010>, 2013.
- Rupp, K.: Microprocessor Trend Data, GitHub [data set], <https://github.com/karlrupp/microprocessor-trend-data> (last access: 12 September 2024), 2022.
- Shchepetkin, A. F. and McWilliams, J. C.: A method for computing horizontal pressure-gradient force in an oceanic model with a nonaligned vertical coordinate, *J. Geophys. Res.*, 108, 3090, <https://doi.org/10.1029/2001JC001047>, 2003.
- Silvestri, S., Wagner, G. L., Constantinou, N. C., Hill, C. N., Campin, J.-M., Souza, A. N., Bishnu, S., Churavy, V., Marshall, J. C., and Ferrari, R.: A GPU-based ocean dynamical core for routine mesoscale-resolving climate simulations, *ESSOAr* [preprint], <https://doi.org/10.22541/essoar.171708158.82342448/v1>, 2024.
- Sridhar, A., Tissaoui, Y., Marras, S., Shen, Z., Kawczynski, C., Byrne, S., Pamnany, K., Waruszewski, M., Gibson, T. H., Kozdon, J. E., Churavy, V., Wilcox, L. C., Giraldo, F. X., and Schneider, T.: Large-eddy simulations with ClimateMachine v0.2.0: a new open-source code for atmospheric simulations on GPUs and CPUs, *Geosci. Model Dev.*, 15, 6259–6284, <https://doi.org/10.5194/gmd-15-6259-2022>, 2022.
- Strohmaier, E., Dongarra, J., Simon, H., Meuer, M., and Meuer, H.: The Top 500, PROMETEUS Professor Meuer Technologieberatung und -Services GmbH, <https://www.top500.org/lists/top500/list/2024/11/> (last access: 14 April 2025), 2024.
- Tsujino, H., Motoi, T., Ishikawa, I., Hirabara, M., Nakano, H., Yamanaka, G., Yasuda, T., and Ishizaki, H.: Reference manual for the Meteorological Research Institute Community Ocean Model (MRI.COM) version 3. Technical Reports of the Meteorological Research Institute, report no. 59, <https://doi.org/10.11483/mritechrepo.59>, 2010.
- Voosen, P.: Climate modelers grapple with their own emissions, *Science*, 384, 494–495, <https://doi.org/10.1126/science.adq1772>, 2024.