# Unlocking the Power of Parallel Computing: GPU technologies for Ocean Forecasting

Andrew Porter[1] and Patrick Heimbach[2]

[1]Science and Technology Facilities Council, Daresbury Laboratory, Hartree Centre, Daresbury, UK
[2]Oden Institute for Computational Engineering and Sciences, The University of Texas at Austin, USA

*Correspondence to*: Andrew Porter (andrew.porter@stfc.ac.uk)

**Abstract.** Operational ocean forecasting systems 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, technological evolution of microprocessors has constrained data parallel scientific computation. Today, GPUs offer an additional and valuable source of computing power to 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/volume/element modelling through to the training of deep neural networks used in machine learning 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, differentiating among tools that perform code-transformation, easing the adaptation of Fortran code for GPU execution (like PSyclone) or direct use of OpenACC directives (like ICON-O), to adoption of specific frameworks that facilitate the management of parallel execution across different architectures.
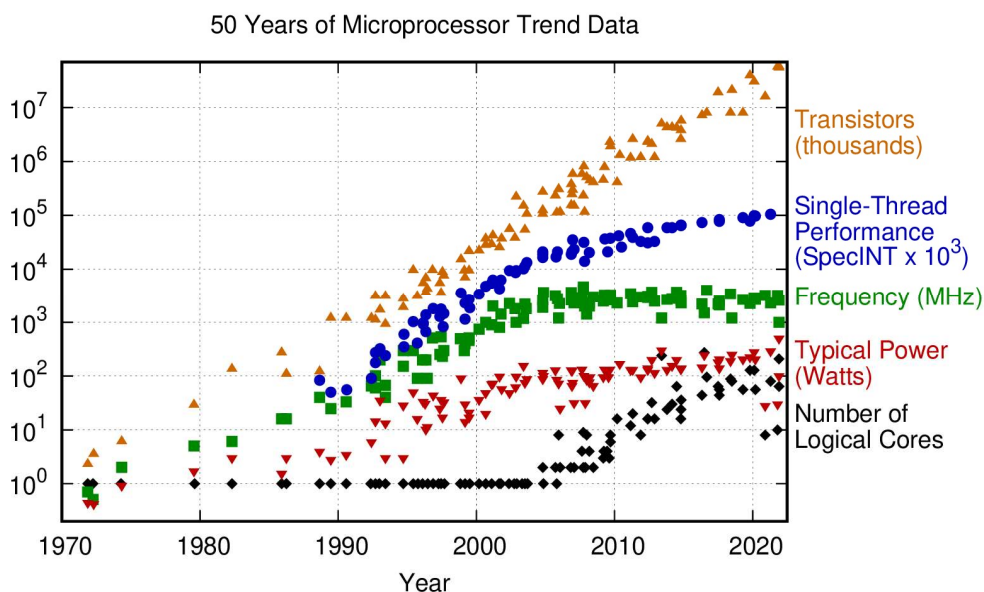
## 1 Introduction

Operational Ocean Forecasting Systems (OOFS) are computationally demanding, and large compute 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 improvement from one generation of processor to the next (Figure 1). Although the number of transistors per device continues to rise, they are increasingly being

STATE OF THE
PLANET
Discussions

31      used to implement larger numbers of execution cores. It is then the job of the application to make use of these additional cores

32      to achieve a performance improvement. Graphical Processing Units (GPUs) are a natural consequence of this evolution.

33      Originally developed to accelerate rendering of computer-generated images (a naturally data-parallel task thanks to the division

34      of an image into pixels), scientists were quick to seize on their potential to accelerate data-parallel scientific computation.

35      Therefore, manufacturers today produce HPC-specific "GPUs" that are purely intended for computation. The suitability of this

36      hardware for the training of deep neural networks used in machine learning and artificial intelligence has stimulated massive

37      development and competition amongst GPU vendors.

38      Unlike CPUs which tend to have relatively few but powerful (general purpose) processor cores, GPUs support hundreds of

39      simpler cores running thousands of threads which can get data from memory very efficiently. The simplicity of these cores

40      makes them more energy efficient and therefore GPUs tend to offer significantly greater performance per Watt. With energy

41      consumption of large computing facilities now the key design criterion, GPUs are an important part of the technology being

42      used in the push towards Exascale performance and beyond (e.g. Draeger and Siegel, 2023). As an illustration, in the November

43      2022 incarnation of the Top500 list (Strohmaier et al., 2020), eight of the machines in the top ten are equipped with some form

44      of accelerator and the majority of those are GPUs from either NVIDIA or AMD. Although CPUs are present in these machines,

45      their primary role is to host the GPUs which provide the bulk of the compute performance. GPUs are therefore a major feature

46      of the current HPC landscape, and their importance and pervasiveness is only set to increase.

Figure 1: 50 years of microprocessor (CPU) evolution showing the breakdown of Dennard scaling (Rupp, 2022)

## 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 are solved numerically by discretizing 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 Same 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.

3

## 3 The use of GPUs in Ocean Forecasting

61 Although GPUs are now a well-established HPC technology with potentially significant performance advantages for OOFS,
62 they are not yet widely adopted in the ocean-forecasting community. For example, in Europe, NEMO (Madec et al., 2023) is
63 the most important ocean-modeling framework; it is used operationally by Mercator Ocean International, the European Centre
64 for Medium-Range Weather Forecasting, the UK Met Office and the Euro-Mediterranean Centre on Climate Change, and
65 other Institutes worldwide. NEMO is implemented in Fortran and parallelised with MPI and as such is limited to running on
66 CPUs only. The German weather service (DWD) uses ICON-Ocean (Korn, 2017) which is also a Fortran model. Experiments
67 are in progress with the use of OpenACC directives to extend this code to make use of GPUs but this functionality is not used
68 operationally.
69 In the US, NOAA's Real-Time Ocean Forecast System (https://polar.ncep.noaa.gov/global/) is based on HYCOM (HYbrid
70 Coordinates Ocean Model, Chassignet et al., 2009). HYCOM too is a Fortran code parallelised using a combination of OpenMP
71 and MPI. Although not used operationally, the Energy Exascale Earth System Model is also significant. It utilizes the MPAS
72 (Model for Prediction Across Scales) Ocean, Sea-Ice and Land-Ice models (Ringler et al., 2013) which again is implemented
73 in Fortran with MPI (although some experimental ports have been performed using OpenACC directives). The MIT General
74 Circulation Model (MITgcm, Marshall et al., 1997) is also widely used and again is Fortran with support for distributed- and
75 shared-memory parallelism on CPU.
76 The Japanese Meteorological Agency runs operational forecasts using the Meteorological Research Institute Community
77 Ocean Model (MRI.COM) (Tsujino et al., 2010). As with the previous models, this too is implemented in Fortran with MPI
78 and thus only runs on CPU.
79 For regional (as opposed to global) forecasts, the Rutgers Regional Ocean Modeling System (ROMS) (Shchepetkin and
80 McWilliams, 2023) is used by centers worldwide including the Japan Fisheries Research and Education Agency, the Australian
81 Bureau of Meteorology and the Irish Marine Institute. ROMS too is a Fortran code parallelised using either MPI or OpenMP
82 (but not both combined) and thus is restricted to CPU execution. Although various projects have ported the code to different
83 architectures (including the Sunway architecture for China's Tianhe machine, Liu et al., 2019), these are all standalone pieces
84 of work that have not made it back into the main code base.

## 4 Discussion

86 From the preceding section, it is clear that OOFS are currently largely implemented in Fortran with no or limited support for
87 execution on GPU devices. The problem here is that OOFS comprise of large and complex codes which typically have a
88 lifetime of decades and are constantly being updated with new science by multiple developers. Maintainability, allowing for
89 the fact that the majority of developers will be specialists in their scientific domain rather than in HPC, is therefore of vital
90 importance. Given that such codes are often shared between organizations, they must also run with good performance on
91 different types of architecture (i.e. be 'performance portable').

92    Previously, one generation of supercomputers looked much like the last and therefore the evolution of these computer models

93    was not a significant problem. However, the proliferation of computer hardware (and, crucially, the programming models

94    needed to target them) that has resulted from the breakdown of Dennard scaling has changed this. With the average

95    supercomputer having a lifetime of just some five years, OOFS are now facing the problem of adapting to future supercomputer

96    architectures and this is difficult because the aims of performance, performance portability and code maintainability often

97    conflict with each other (Lawrence et al., 2018).

98    To date there have been various approaches to this problem. NEMO is in the process of adopting the PSyclone code-

99    transformation tool (Adams et al., 2019) that enables an HPC expert to transform Fortran source code such that it may be

100   executed on GPU using whichever programming model is required. For a low-resolution, 1 degree) global mesh, a single

101   NVIDIA V100 GPU gives a performance some 3.6x better than an HPC-class Intel socket. For a high-resolution, (1/12th

102   degree) global mesh, ~90 A100 GPUs give the same performance as ~270 Intel sockets (Porter et al., 2023 - in prep.). As noted

103   earlier, ICON-O is being extended manually with OpenACC directives (although these are only supported on NVIDIA

104   hardware). There are examples of recent (i.e. experimental) models that have moved away from Fortran in favor of higher-

105   level programming approaches. Thetis (Kärnä et al., 2018) implements a Discontinuous Galerkin method for solving the 3D

106   hydrostatic equations using the Firedrake framework. This permits the scientist to express their scheme in the Python

107   implementation of Unified Form Language (Alnæs et al., 2014). The necessary code is then generated automatically. The

108   Veros model (Häfner et al., 2021) takes a slightly different approach: its dynamical core is a direct Python translation of a

109   Fortran code and thus retains explicit MPI parallelisation. The JAX system (http://github.com/google/jax) for Python is then

110   used to generate performant code for both CPU and GPU. The authors report that the Python version running on 16 A100

111   GPUs gives the same performance as 2000 CPU cores for the Fortran version (although this comparison is slightly unfair as

112   the CPUs used are several generations older than the GPUs).

113   Another popular approach to performance portability is to implement a model using a framework that takes care of parallel

114   execution on a target platform. Frameworks such as Kokkos (Carter Edwards et al., 2014), SyCL and OpenMP are good

115   examples. In principle this approach retains single-source science code, while enabling portability to a variety of different

116   hardware. However, it is hard to insulate the oceanographer from the syntax of the framework (which are often only available

117   in C++) and, while the framework may be portable, obtaining good performance often requires that it be used in a different

118   way from one platform to another. In OpenMP for instance, the directives needed to parallelise a code for a multi-core CPU

119   are not the same as those needed to offload code to an accelerator.

120   The Climate Modeling Alliance (CliMA) has adopted a radically new approach by rewriting ocean and atmospheric models

121   from scratch using the programming language Julia (Perkel, 2019; Sridhar et al., 2022). Designed to overcome the "two-

122   language problem" (Churavy et al., 2022), Julia is ideally suited to harness emerging HPC architectures based on GPUs (Besard

123   et al., 2017; Bezanson et al., 2017). First results with CliMA's ocean model, Oceananigans.jl (Ramadhan et al., 2020), run on

124   64 NVIDIA A100 GPUs exhibit 10 Simulation Years Per Day (SYPD) at 8 km horizontal resolution (Silverstri et al., 2024).

125   This performance is similar to current-generation CPU-based ocean climate models run at much coarser resolution (order 25-

126    50 km resolution). Similarly promising benchmarks have been obtained with a barotropic configuration of a prototype of the

127    MPAS-Ocean model, rewritten in Julia (Strauss et al., 2023). Such performance gains hold great promise for accelerating

128    operational ocean prediction at high spatial resolution run on emerging HPC hardware.

129    **References**

130    Adams, S. V., Ford, R. W., Hambley, M., Hobson, J. M., Kavcic, I., Maynard, C. M., Melvin, T., Mueller, E. H., Mullerworth,

131    S., Porter, A. R., Rezny, M., Shipway, B. J. and Wong, R.: LFRic: Meeting the challenges of scalability and performance

132    portability in Weather and Climate models. Journal of Parallel and Distributed Computing, 132, 383-396.

133    https://doi.org/10.1016/j.jpdc.2019.02.007, 2019.

134    Alnæs, M. S., Logg, A., Ølgaard, K. B., Rognes, M. E. and Wells, G. N.: Unified Form Language: A Domain-Specific

135    Language for Weak Formulations of Partial Differential Equations, ACM Trans. Math. Softw., 40(2), doi:10.1145/2566630

136    Carter Edwards, H., Trott, C.R., and Sunderland, D. (2014). Kokkos: enabling manycore performance portability through

137    polymorphic memory access patterns, Journal of Parallel and Distributed Computing, 74(12), 3202-3216.

138    https://www.osti.gov/servlets/purl/1106586, 2014.

139    Besard, T., Foket, C. & Sutter, B. D.: Effective Extensible Programming: Unleashing Julia on GPUs. IEEE Transactions on

140    Parallel and Distributed Systems, 30(4), 827–841. https://doi.org/10.1109/tpds.2018.2872064, 2017.

141    Bezanson, J., Edelman, A., Karpinski, S. and Shah, V. B.: Julia: A Fresh Approach to Numerical Computing. SIAM Review,

142    59(1), 65–98. https://doi.org/10.1137/141000671, 2017;

143    Chassignet, E.P., Hurlburt, H.E., Metzger, E.J., Smedstad, O.M., Cummings, J., Halliwell, G.R., Bleck, R., Baraille, R.,

144    Wallcraft, A.J., Lozano, C., Tolman, H.L., Srinivasan, A., Hankin, S., Cornillon, P., Weisberg, R., Barth, A., He, R., Werner,

145    F., and Wilkin, J.: U.S. GODAE: Global Ocean Prediction with the HYbrid Coordinate Ocean Model (HYCOM).

146    Oceanography, 22(2), 64-75. https://doi.org/10.5670/oceanog.2009.39, 2009.

147    Churavy, V., Godoy, W. F., Bauer, C., Ranocha, H., Schlottke-Lakemper, M., Räss, L., Blaschke, J., Giordano, M., Schnetter,

148    E., Omlin, S., Vetter, J. S. & Edelman, A.: Bridging HPC Communities through the Julia Programming Language. arXiv.

149    https://doi.org/10.48550/arxiv.2211.02740, 2022.

150    Dennard, R. H., Gaensslen, F., Yu, H., Rideout, L., Bassous, E., and LeBlanc, A.: Design of ion-implanted MOSFET's with

151    very small physical dimensions. IEEE Journal of Solid-State Circuits. SC-9 (5), 256–268, 1974.

152    Draeger, E. W. and Siegel, A.: Exascale Was Not Inevitable; Neither Is What Comes Next. Computing in Science and

153    Engineering, 25(3), 79–83. https://doi.org/10.1109/mcse.2023.3311375, 2023.

154    Häfner, D., Nuterman, R. and Jochum, M.: Fast, cheap, and turbulent — Global ocean modeling with GPU acceleration in

155    Python. Journal of Advances in Modeling Earth Systems, 13. doi:10.1029/2021MS002717, 2021.

STATE OF THE
PLANET
Discussions

156   Kärnä, T., Kramer, S. C., Mitchell, L., Ham, D. A., Piggott, M. D. and Baptista, A. M. .: Thetis coastal ocean model:
157   discontinuous Galerkin discretization for the three-dimensional hydrostatic equations. Geoscientific Model Development,
158   11(11), 4359-4382. https://doi.org/10.5194/gmd-11-4359-2018, 2018.

159   Korn, P.: Formulation of an unstructured grid model for global ocean dynamics, Journal of Computational Physics, 339, 525-
160   552. https://doi.org/10.1016/j.jcp.2017.03.009, 2017.

161   Lawrence, B. N., Rezny, M., Budich, R., Bauer, P., Behrens, J., Carter, M., Deconinck, W., Ford, R., Maynard, C.,
162   Mullerworth, S., Osuna, C., Porter, A., Serradell, K., Valcke, S., Wedi, N. and Wilson, S.: Crossing the chasm: how to develop
163   weather and climate models for next generation computers? Geoscientific Model Development, 11(5), 1799-1821.
164   https://doi.org/10.5194/gmd-11-1799-2018, 2018.

165   Madec, G. and NEMO System Team: NEMO ocean engine, Scientific Notes of Climate Modelling Center, 27, Institut Pierre-
166   Simon Laplace (IPSL), doi:10.5281/zenodo.1464816, 2023.

167   Marshall, J., Adcroft, A., Hill, C., Perelman, L., and Heisey, C.: A finite-volume, incompressible Navier Stokes model for
168   studies of the ocean on parallel computers. J. Geophysical Res., 102(C3), 5753-5766. https://doi.org/10.1029/96JC02775,
169   1997.

170   Perkel, J. M.; Julia: come for the syntax, stay for the speed. 141–142. http://www.nature.com/articles/d41586-019-02310-3,
171   2019.

172   Porter, A. R., Dearden, C., Ford, R. W., Henrichs, J., Siso, S., Nobre, N., Müller, S. A., Coward, A., Bell, M. (2023). Using
173   PSyclone 3.4 to achieve Performance Portability for NEMO 4.0 and NEMO-MEDUSA Ocean Models, in preparation.

174   Ringler, T. Petersen, M., Higdon, R. L., Jacobsen, D., Jones, P. W. and Maltrud, M.: A multi-resolution approach to global
175   ocean modeling, Ocean Modelling, 69, 211-232. https://doi.org/10.1016/j.ocemod.2013.04.010, 2013.

176   Shchepetkin, A. F., and McWilliams, J.C.: A method for computing horizontal pressure-gradient force in an oceanic model
177   with a nonaligned vertical coordinate. Journal of Geophysical Research: Oceans, 108(C3), 3090.
178   https://doi.org/10.1029/2001JC001047, 2003.

179   Ramadhan, A., Wagner, G., Hill, C., Campin, J.-M., Churavy, V., Besard, T., Souza, A., Edelman, A., Ferrari, R. & Marshall,
180   J.: Oceananigans.jl: Fast and friendly geophysical fluid dynamics on GPUs. Journal of Open Source Software, 5(53), 2018.
181   https://doi.org/10.21105/joss.02018, 2020.

182   Silvestri, S., Wagner, G. L., Constantinou, N. C., Hill, C. N., Campin, J.-M., Souza, A. N., Bishnu, S., Churavy, V., Marshall,
183   J. C. and Ferrari, R.: A GPU-based ocean dynamical core for routine mesoscale-resolving climate simulations. ESSOAr.
184   https://doi.org/10.22541/essoar.171708158.82342448/v1, 2024.

185   Sridhar, A., Tissaoui, Y., Marras, S., Shen, Z., Kawczynski, C., Byrne, S., Pamnany, K., Waruszewski, M., Gibson, T. H.,
186   Kozdon, J. E., Churavy, V., Wilcox, L. C., Giraldo, F. X. & Schneider, T.: Large-eddy simulations with ClimateMachine
187   v0.2.0: a new open-source code for atmospheric simulations on GPUs and CPUs. Geoscientific Model Development, 15(15),
188   6259–6284. https://doi.org/10.5194/gmd-15-6259-2022, 2022.

STATE OF THE
PLANET
Discussions

189 Strauss, R. R., Bishnu, S. and Petersen, M. R.: Comparing the Performance of Julia on CPUs versus GPUs and Julia-MPI
190 versus Fortran-MPI: a case study with MPAS-Ocean (Version 7.1). EGUsphere, 2023, 1–22.
191 https://doi.org/10.5194/egusphere-2023-57, 2023.
192 Strohmaier, E., Dongarra, J., Simon, H., Meuer, M. and Meuer, H.: The Top 500.
193 https://www.top500.org/lists/top500/list/2022/11/, 2022.
194 Tsujino, H., Motoi, T., Ishikawa, I., Hirabara, M., Nakano, H., Yamanaka, G., Yasuda, T., and Ishizaki, H.: Reference manual
195 for the Meteorological Research Institute Community Ocean Model (MRI.COM) version 3. Technical Reports of the
196 Meteorological Research Institute, 59, 273, 2010.
197 Rupp, K: Microprocessor Trend Data, https://github.com/karlrupp/microprocessor-trend-data, 2022. Accessed 12/09/2024.

198 **Competing interests**

199 Author A. Porter has declared that he is an author of the PSyclone package and a co-chair of the NEMO HPC Working Group.

200 **Data and/or code availability**

201 This can also be included at a later stage, so no problem to define it for the first submission.

202 **Authors contribution**

203 This can also be included at a later stage, so no problem to define it for the first submission.

204 **Acknowledgements**

205 This can also be included at a later stage, so no problem to define it for the first submission.