**To editor (minor revision)**

1. **The authors have done a good job addressing most of the issues raised by the referees, and I largely agree with their judgement in choosing to not make some suggested revisions. The only case that I have some disagreement with is their decision to not make any revisions to address the anonymous referee's comments about the superlinearity discussion in Section 6.1. I agree with the referee that this discussion would be strengthened by providing some additional measurements collected with hardware counters to demonstrate the authors' belief that the superlinear speedup is due to improved cache utilization. However, I do not think this is critical, and I respect the authors' decision to not do so. However, I think that the wording here should be amended to make it clear that the authors believe that the superlinear speedup is due to cache effects, but that this has not been verified. Where the authors state "We explain superlinearity by better cache usage", the wording may lead some readers to believe that this explanation has been verified either experimentally or using some cache model. I suggest changing this to something like "We believe the superlinear speedup is due to improved cache usage, but we have not investigated this". It may also be useful to incorporate some of the content they wrote in their justification of their response to the referee on this point, though I leave this to the authors' discretion.**

We greatly appreciate the editor's and anonymous referee's effort to clarify this point. It was helpful to investigate on-node performance using Intel Advisor profiler. Superlinearity occurs in the range of cores 1-149. We have found that for loops in the function iceadvect the time of waiting data from the memory is distributed in the following way for low number of cores:



and for 149 cores:



I.e., memory delay is limited by the DRAM latency and\or bandwidth on low number of cores, and starting from 149 cores most memory accesses correspond to the L1 cache. Analogous behavior for the advect3D function is less crucial, since it is more computationally expansive and rely more on the L1 cache. We think that strict quantitative analysis of data transfer should imply some subsequent reorganization of the calculations, and should be addressed in the future, as the anonymous referee suggests. Thus, it seems enough to give a short qualitative description of this point, and, so, we add the following discussion in Section 6.1.:

"Superlinearity (parallel efficiency is greater than 100% when "doubling" the number of cores) occurs in the range of cores 1-149, and we have investigated the possible reasons for this using the open access Intel Advisor profiler. We have found that on 1 core the most time demanding memory requests are the DRAM data upload, while on 149 cores most memory requests correspond to the L1 cache. Thus, superlinearity can be partially explained by the better cache utilization when the number of cores increases. Note that more local memory access pattern (MAP) can decrease the limitations caused by the memory requests and can be achieved by incorporating stride-1 access for the inner loop indices, but we leave this point for optimizations in the future."

Also, for readability, this discussion is continued with a new paragraph which starts as:

"Speedup including MPI exchanges is shown in figure 4 with solid lines."

**To editor**

Dear editor,

We greatly appreciate your comment. Version (FEMAO 2.0) is presented in the revised manuscript.

Best regards, On behalf of co-authors, Pavel Perezhogin

**To Koldunov**

We greatly thank the reviewer for the careful reading of this manuscript and given suggestions.

1. **The two step procedure of first dividing the model domain into small blocks and then redistributing those blocks between cores was not really clear to me at first and should be better communicated. It would be helpful for uninitiated readers if you can mention earlier on that the requirement is to preserve the structured nature of the code. So your partitions can't be of arbitrary shape, like in unstructured mesh models, but should be constructed out of small rectangles. I would suggest creating a schematic that shows all the steps of the procedure - splitting into so called blocks, fitting the Hilbert curve, distributing the blocks among CPU cores and finally allocating "shared" arrays. Of course it's not possible to demonstrate with 128x128 blocks you use for a realistic model, but something like a 10x10 schematic representation would do the job.**
   **A bit more details on how the partitioning handled in the model setup would be appreciated. Does the partitioning created by the library and then read by the model? Or it's computed each time. If the latter is the case - do you guarantee that the partitioning will be the same each time the model is run?**

   We do not think that there is a need to additionally explain algorithm of distribution of the blocks over the cores, because it doesn't meet the main objective of the paper, have been shown many times by Dennis and there is a general-purpose solution (METIS). "Shared" arrays are clarified in figures 2 and 3. There is no need for another figure.

   The introduction is changed:
   P2 L34 "In numerical ocean model…" is moved to new paragraph
   P2 L39 "We give preference …" is removed
   P2 L41 "Note that some modern…" is moved to previous paragraph

   We add the last paragraph in the introduction:
   "In sections 2-4 we provide model configuration and organization of the calculations in the non-parallel code on structured rectangular grid. In section 5 we describe parallelization approach, which preserves original structure of the loops. Domain decomposition is carried out in two steps: first the model domain is divided into small blocks and then these blocks are distributed between CPU cores. For all blocks belonging to a given core a "shared" array is introduced, and mask of computational points restricts calculations. Partition could be of arbitrary shape, but blocks allow us to reach the following benefits: simple balancing algorithm (Hilbert curves) can be applied as the number of blocks along a given direction is chosen to be a power of 2; boundary exchanges can be easily constructed for arbitrary halo width, but smaller than the block size. In section 6 we report parallel acceleration on different partitions for particular 2D and 3D subroutines and the whole model."

Section name "Organization of the calculations" is changed to "Organization of the calculations in non-parallel code".

We add the first paragraph to the section "Modifications of the non-parallel code":
"In this section we describe the partitioning algorithm of the model domain into subdomains, each corresponding to a CPU core, and subsequent modifications of the single-core calculations, which require only minor changes of algorithms 2 and 3. Grid partition is performed in two steps: model domain is decomposed into small blocks and then these blocks are distributed over CPU cores in such a way that computational load imbalance is minimized. We utilize common grid partition for both sea-ice and ocean submodels, and provide theoretical estimates of the load imbalances resulting from the application of different weight functions in the balancing problem. Partition is calculated during the model initialization stage, as our balancing algorithm (Hilbert curves) is computationally unexpensive. Also, we guarantee that the partition is the same each time the model is run, if parameters of the partitioner were not modified. "

2. **Minor comments.**
   1) Line 23: "adjusted to the White Sea Chernov (2013), Chernov et al. (2018)" - You forgot parentesis.
      **Response**: We add parenthesis.
   2) Line 28:"(i.e., not sigma coordinate and so on)." - Just delete it, you don't need this clarification.
      **Response**: Deleted.
   3) Line 28:"In case of significantly variable depth, this "integer depth" may also vary, see figure 1." - I think I understand what you are trying to say here, the number of levels vary with depth, but it's not clear why depth should be "significantly variable". Please rewrite to make it clearer.
      **Response**: Rewritten: "In case of significantly variable depth, the number of levels also varies, see figure 1."
   4) Line 30 - what balancing (I assume you mean balancing of model computation)?
      **Response:** Rewritten: "The presence of both 2D and 3D calculations complicates *balancing of the computations* for the full model."
   5) Line 39: "distributed using the METIS Karypis (1998)" - you need parentheses around citation in this case. Please double check all your citations.
      **Response:** parentheses are added. We have checked all citations.
   6) Line 40: "to make the code library-independent" - my understanding is that you create a separate library for partitioning, so at the end it depends on the library, it's just your library? :)
      **Response:** This sentence is removed. Instead, we add the final paragraph in the Introduction with more accurate description of our approach.
   7) Line 57:"initial distributions" - of what? Please be more precise.
      **Response:** Sentence is rewritten: "Second, because this model is less dependent on the initial data, it makes the test simulations easier because the only liquid boundary is needed to set the initial-boundary data."
   8) Line 59:"demonstrate any important features" - please rephrase, maybe give some examples.
      **Response:** Rewritten: "Finally, the White Sea's relatively small inertia enables to check correctness of the code by rather short simulations, which are able to demonstrate important features of the currents."
   9) Line 66 - Please provide details on the type of advection you use.

**Response:** Line 73 of the revised manuscript: "The simple Characteristic-Galerkin Scheme (Zienkiewicz and Taylor,2000) is used for the 3D and 2D advection terms."

10) Lines 68-73. - Please provide references for sea ice dynamics and thermodynamics.

**Response:** Line 75 of the revised manuscript. Paragraph is rewritten:

"The local 1D sea ice thermodynamics is based on the 0-layer model (Semtner, 1976; Parkinson and Washington, 1979) with some modifications in lateral melting and surface albedo (Yakovlev, 2009). There are 14 categories of ice thickness (gradations), the mechanical redistribution and the ice strength are identical to the CICE (Hunke et al., 2013). The elastic-viscous plastic scheme (EVP; Danilov et al., 2015) with modification for the relaxation time scales (Wang et al., 2016) is used for the sea ice dynamics (see also the Appendix 3 in Koldunov et al., 2019b). Sea ice is described by distribution of 80 its compactness (concentration) and ice volume for each gradation. In addition, snow-on-ice volume for each gradation is evaluated. Therefore, there are 43 2D sea-ice scalars: ice and snow volume for 14 gradations and sea-ice compactness for 15 ones (including water). Because there are 39 vertical layers in an ocean component, the set of all of the sea-ice data is comparable to a single 3D scalar."

11) Line 77:"more shallow than it really is" - any references to that?

**Response:** Reference is provided: "Comparison of available bathymetry data for the White Sea is given in Chernov and Tolstikov (2020) in table 1."

12) Line 102 - What do you mean by subdomain? Number of blocks that belong to one core? subdomain in computational domain, like a bay? Please define.

**Response:** Rewritten: "Connectivity of subdomains (by subdomain we refer to a set of blocks belonging to a CPU core) or minimum length of the boundary can be chosen as possible criteria for the quality of a partition."

13) Line 129:"Let us introduce two baseline partitions:" - change to "We have implemented two baseline partitions"

**Response:** changed.

14) Fig. 2:"Black rectangle corresponds to a "shared" array" - change to "Black rectangle on figure for hilbert3d partition corresponds to ..."

**Response:** Rewritten: "Black rectangle in figure for hilbert3d partition corresponds to …"

15) Line 151:"The shared array size is shown..." - Change to "An example of the shared array size..."

**Response:** changed.

16) Line 189:"Simulations were performed for three model days" ** and then in Line 192 "launched on 993 CPU cores for 30 days" - Please clarify.

**Response:** These sentences are rewritten:

"Low-core simulations were performed for three model days (2592 time steps). The model on 993 CPU cores is launched for 30 days, with subsequent rescaling of the results."

**To anonymous referee.**

We are grateful to the referee for the very helpful comments and given suggestions.

1. **The main subject of the paper is the MPI implementation and load balancing, but I suspect some aspects of the model are limiting on-node performance. For example, they describe their choice of conditional masking vs multiplicative masking for land points (pg 2 line 42), their non-optimal combination of loop and index ordering (pg 2 L90), and the potential advantages of unstructured meshes (p2, L41). They have included at least some discussion of these in the paper so I'm not suggesting any changes now, but may have some implications on later comments below and would encourage them to explore these as they continue their optimizations in the future.**

In this model, boundary conditions are included into matrix elements, which are stored as an array KT(6,13), where the first dimension corresponds to 6 triangles composing Finite Element, and the second dimension corresponds to 13 types of "wet" points: 1 inside the domain and 12 types of boundary points. This approach is similar to multiplicative masking, as B.C.s are applied by the product to KT, and to unstructured mesh models, as matrix elements are precomputed. The difference from the unstructured mesh models is that only unique elements of the matrix are stored and neighboring points are referenced directly. So, the mask of wet points serves only to restrict the number of computations. We thank the referee and will think in the future how to organize calculations more efficiently.

In our opinion, for the model we have for now, it is not reasonable to change loops' order, as it harms model infrastructure and our parallelization approach, but array indices may be chosen more optimally, setting "depth" index as the first. Nevertheless, this interchange is not crucial for the goals we address in the paper.
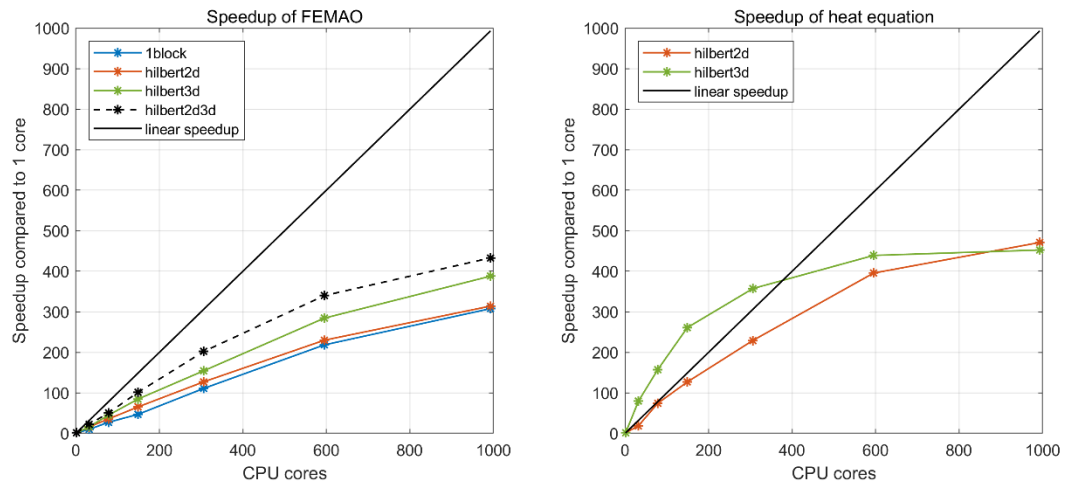
**Response:** without changes.

2. **First, the authors show speedups in figure 4 with significantly super-linear speedups in the 2-d case. They attribute this to cache performance without additional evidence (eg from hardware counters or other performance tools). That may be the case, but I think this super-linearity is large enough to warrant further exploration into the cause.**

This super-linear speedup is measured for the code section of approximate length 1000 code lines which consist of 6 loops like in algorithm 3. We stress that these loops have slightly different organization of the calculation and may accelerate in slightly different rates. Some of the loops work with 4D arrays, where the first additional dimension corresponds to "antidiffusive fluxes". Some loops have additional if-conditions, which are needed to perform flux correction in quasi-monotone scheme. Superlinearity occurs at the low-to-middle number of cores, and these cases are usually omitted when scaling up to many cores is shown. Moreover, usually speedups are shown including exchanges, and for this option our speedups are not superlinear.

Finally, from the practical viewpoint, the presented parallelization approach together with the chosen loops/indices ordering may lead to superlinear acceleration. As an example, consider very simple "heat equation" loop:

```
do j = js, je
    do i = is, ie
        do k = 1,depth(i,j)
            Tn(i,j,k) = 0.25_8 * (T(i+1,j,k) + T(i-1,j,k) + T(i,j+1,k) + T(i,j-1,k))
        end do
    end do
end do
```

It accelerates superlinearly at the low-to-middle number of cores for appropriate weights even when MPI exchanges are taken into account (green line in the right subfigure):



Regardless of the actual reason it happens (decrease in the number of cache misses, non-optimal organization of the calculations, or something else), this "heat equation" loop constitutes what we actually intended to do, and there is nothing to optimize here.

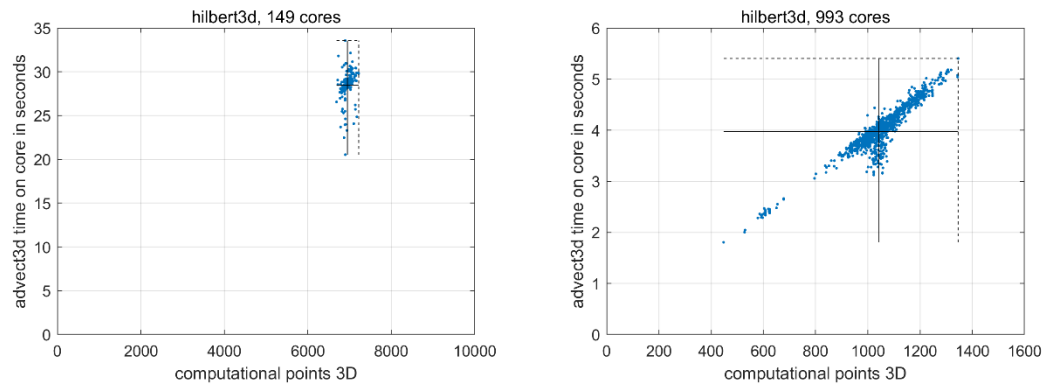**Response:** without changes (see also response to editor).

3. **Third, the large variation in work load at high core counts (fig 5,6) also seems higher than one might expect. As you get fewer points/blocks per core, there will naturally be a little higher variability, but this seems larger than expected and might point to additional problems.**

As the referee pointed out, balancing could be better. In experiment presented in the manuscript, balancing is limited by the outlier point (figs. 5 and 6), corresponding to the CPU core with 5 blocks and maximum load. This outlier point limits LI to 46%. We have checked balancing optimization procedure (algorithm 4) and found that it doesn't guarantee monotone decrease of LI, as subroutine "remove_not_connected_subdomain" can increase LI. After choosing the best iteration, LI was decreased to 29%. As this behavior is crucial only for "hilbert3d 993 cores" experiment, in the revised manuscript we will update it. Additionally, we have tested METIS multilevel k-way contiguous partitioning algorithm and found that it doesn't give better balancing (LI=39%).

**Response:** "hilbert3d 993 cores" experiment, which is presented in table 3 and figures 4-9, was updated.

4. **Second, the computational time as a function of wet points seems a bit counter-intuitive (Fig. 5). The authors have shown percentage of wet points rather than total wet points to emphasize their diagnosis again of memory access. But without also seeing the total number of points (computational load), it's a little hard to get a more complete picture. Again, this effect seems too big to attribute solely to cache effects and it seems like more might be going on here.**

Figures 5 and 6 are provided to assess separately data structure efficiency and load balancing efficiency and clearly show limitations of the described model. As the referee pointed out, the lack of point-to-point correspondence between 5 and 6 figures lead to incomplete picture of what is going on. Here we provide scatter plot (6 figure y axis – 5 figure y axis):

Scatterplots are provided with mean values (solid lines) and maximum values (dashed lines). These values completely define Load Imbalance (LI) in partition and advect3d runtime. As follows from the left figure, spread in runtime is more then spread in the number of computational points. This means that computations are limited by the organization of the calculations, but not by the accuracy of the partitioning algorithm. As advect3d is a function with approximate length of 2500 code lines, which consists of 6 loops like in algorithm 2, and each loop has slightly different organization of the calculations, we claim that overestimation of runtime LI only by 15% in comparison to partition LI is a very good result. In the right figure, there is strict correlation between the number of computational points and advect3d runtime, and computations are limited by the balancing procedure. As this figure is more informative than figure 6, in the revised manuscript we will attach the new figure.

**Response:** Figure 6 is changed to the new one. Also, we update discussion of this figure (lines 247-254 in the revised manuscript):

"Figure 6 additionally shows that spread in runtime cannot be explained by the difference in the number of computational points, i.e. partitioning algorithm works well for 149 cores. Although organization of the calculations may slightly limit model efficiency on a moderate number of cores, it does not limit the model efficiency on 993 cores, where major part of the advect3D runtime spread is explained by the imperfect balancing (see figure 6), but not the data structure (see figure 5). Stagnation of the balancing procedure is evident from the fact that the minimum number of blocks located on a CPU core is 1 for 993 cores, see table 3. Note that computational subdomain corresponding to one CPU core is small enough: on average, it has 9x9 horizontal points with 12 vertical levels for 993 cores."

5. **Minor comments.**
   1) The journal editor will probably mention this, but most references should be changed so that the parentheses are around both author and date unless an integral part of sentence. So for example p1L22-23, should have (FEMAO; Iakolev, 1996, 2012) and (Chernov, 2013; Chernov et al. 2018). And so on throughout the manuscript.
   **Response:** All references are checked.
   2) Fig 2 has cropped the bottom of figures
   **Response:** Fig 2 is shown as we expected. We do not provide axis labels as they correspond to the mesh points, but not to geographical coordinates.
   3) P9L170-180; In this bulleted list, move the text "These three bullets. . ." and "This reduces. . ." after the bulleted list as "The first three bullets. . ." and "The final bullet. . ." Mixing these comments in with the bulleted list was confusing.
   **Response:** done.

# Advanced parallel implementation of the coupled ocean-ice model FEMAO (version 2.0) with load balancing

Pavel Perezhogin[1], Ilya Chernov[2], and Nikolay Iakovlev[1]

[1]Marchuk Institute of Numerical Mathematics of the Russian Academy of Sciences, Moscow, Russia
[2]Institute of Applied Math Research, Karelian Research Centre of RAS, Petrozavodsk, Russia

**Correspondence:** Pavel Perezhogin, pperezhogin@gmail.com

**Abstract.** In this paper, we present a parallel version of the finite element model of the Arctic Ocean (FEMAO) configured for the White sea and based on the MPI technology. This model consists of two main parts: an ocean dynamics model and a surface ice dynamics model. These parts are very different in terms of the amount of computations because the complexity of the ocean part depends on the bottom depth, while that of the sea-ice component does not. In the first step, we decided to locate both submodels on the same CPU cores with the common horizontal partition of the computational domain. The model domain is divided into small blocks, which are distributed over the CPU cores using Hilbert-curve balancing. Partition of the model domain is static (i.e., computed during the initialization stage). There are three baseline options: single block per core, balancing of 2D computations and balancing of 3D computations. After showing parallel acceleration for particular ocean and ice procedures, we construct the common partition, which minimizes joint imbalance in both submodels. Our novelty is using arrays shared by all blocks that belong to a CPU core instead of allocating separate arrays for each block, as is usually done. Computations on a CPU core are restricted by the masks of not-land grid nodes and block-core correspondence. This approach allows us to implement parallel computations into the model that are as simple as when the usual decomposition to squares is used, though with advances of load balancing. We provide parallel acceleration of up to 996 cores for the model with resolution $500 \times 500 \times 39$ in the ocean component and 43 sea-ice scalars, and we carry out detailed analysis of different partitions on the model runtime.

## 1 Introduction

The increasing performance and availability of multiprocessor computing devices makes it possible to simulate complex natural systems with high resolution, while taking into account important phenomena and coupling comprehensive models of various subsystems. In particular, more precise, accurate, and full numerical description of processes in seas and oceans have become possible. There are now models of seas that can simulate currents, dynamics of thermohaline fields, sea ice, pelagic ecology, benthic processes, and so on; see, for example, review ~~Fox-Kemper et al. (2019)~~ Fox-Kemper et al. (2019).

The finite-element model of the Arctic Ocean ~~(FEMAO) Iakovlev (1996, 2012)~~ (FEMAO; Iakovlev, 1996, 2012) has been developed since the 1990s and it has been adjusted to the White Sea ~~Chernov (2013), Chernov et al. (2018)~~(Chernov, 2013; Chernov et al., . The model domain is a part of the cylinder over sphere (i.e., the Cartesian product of a region on the Earth surface to a vertical

**1**

25   segment). The coordinates are orthogonal, with the axes directed to the East, to the South, and downwards. The horizontal grid is structured and rectangular because finite elements are defined on triangles composing rectangles, see ~~Iakovlev (1996)~~ Iakovlev (1996). Points that correspond to the land are excluded from the computations using a *mask of "wet" points*. The z-coordinate is used as the vertical axis~~(i. e., not sigma coordinate and so on).~~. Therefore, for each 2D-grid node, there is the number of actually used vertical layers. In case of significantly variable depth, ~~this "integer depth" may also vary~~the number of

30   levels also varies, see figure 1. In contrast, sea ice and sea surface computations are depth-independent. The presence of both 2D and 3D calculations complicates balancing of the computations for the full model.

    The original code was written in Fortran-90/95 and it did not allow computation in parallel. Our goal is to develop a parallel version of the model based on the MPI technology without the need to make significant changes in the program code (i.e., preserve loops structure, mask of wet points, but benefit from load balancing). Consequently, we developed a library that

35   performs a partition of the 2D computational domain and organizes communication between the CPU cores.

    In numerical ocean models, the baseline strategy is to decompose domain into squares ~~Madec et al. (2015)~~ (Madec et al., 2015) or into small blocks, with consequent distribution over the processor cores ~~Dennis (2007, 2003); Chaplygin et al. (2019)~~(Dennis, 2007, 200. Both approaches allow to preserve the original structure of the loops and utilize the direct referencing of neighbouring grid nodes on rectangular grids. Decomposition into small blocks is more attractive from the viewpoint of load balancing, especially

40   for z-coordinate models. Blocks can be distributed using the METIS ~~Karypis (1998)~~ (Karypis, 1998) software or simpler algorithms, such as Hilbert curves ~~Dennis (2007). We give preference to partition on blocks, which are distributed using Hilbert curves to make the code library-independent.~~

    (Dennis, 2007). Note that some modern ocean models can also benefit from unstructured mesh usage, where there is no need for the mask of wet points; see for example Koldunov et al. (2019a). In addition, some ocean models omit masking of

45   wet points, see Madec et al. (2015). This implies increase in the number of computations, but benefits from less control-flow interruptions that give rise to better automatic vectorization of loops. ~~In the following sections we will describe our parallel version implementation *relying* on the use of mask of wet points to make balanced computations and we will also outline its peculiar properties~~

    In sections 2-4 we provide model configuration and organization of the calculations in the non-parallel code on structured

50   rectangular grid. In section 5 we describe parallelization approach, which preserves original structure of the loops. Domain decomposition is carried out in two steps: first the model domain is divided into small blocks and then these blocks are distributed between CPU cores. For all blocks belonging to a given core a "shared" array is introduced, and mask of computational points restricts calculations. Partition could be of arbitrary shape, but blocks allow us to reach the following benefits: simple balancing algorithm (Hilbert curves) can be applied as the number of blocks along a given direction is chosen to be a power

55   of 2; boundary exchanges can be easily constructed for arbitrary halo width, but smaller than the block size. In section 6 we report parallel acceleration on different partitions for particular 2D and 3D subroutines and the whole model.
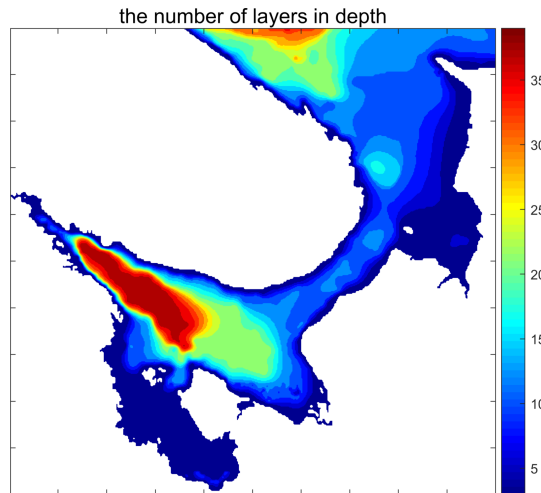
**Figure 1.** The number of depth layers in the White Sea model; the vertical grid has step 5 m up to the 150 m deep and then 10 m up to the 240 m.

## 2 The White Sea

The White Sea is a relatively small (about $500 \text{ km} \times 500 \text{ km}$) and shallow (67 m is the mean depth with a maximal depth of not more than 340 m) semi-closed sea in the Arctic Ocean basin, located in the North-Western part of Russia and included in its territorial waters. Its area is 90000 km$^2$. The White Sea plays an important role for economy of the neighbouring regions ~~Filatov et al. (2007)~~(Filatov et al., 2007).

The White sea consists of several parts, including four bays and a narrow shallow strait called Gorlo that separates one part of the sea from the other. The coastline of the sea is quite complex, which means that the rectangle (almost a square) of the Earth's surface that contains the sea has only about one third of the water area.

The White Sea is a convenient model region to test the numerical algorithms, software, and mathematical models that are intended to be used for the Arctic Ocean. First, low spatial step and relatively high maximum velocities demand, due to the Courant stability condition, a rather small temporal step. This makes it difficult to develop efficient algorithms, stable numerical schemes, and ensure performance using the available computers. Second, because this model is less dependent on the initial ~~distributions~~data, it makes the test simulations easier because the only liquid boundary is needed to set the ~~boundary~~ initial-boundary data. Finally, the White Sea's relatively small inertia enables ~~quite short simulationsthat~~ to check correctness of the code by rather short simulations, which are able to ~~reveal any problems and demonstrate any important features~~ demonstrate important features of the currents.

3

## 3   The model and the software

A time step in the FEMAO model consists of several procedures, see algorithm 1. The model uses the physical-process splitting
75  approach, so that geophysical fields are changed by each procedure that simulates one of the geophysical processes.

---

**Algorithm 1** Time step algorithm for FEMAO

---

1: Forcing (i.e., preparation of river runoff, atmospheric data, shortwave radiation, boundary values, etc.);

2: Dynamics of the sea ice, including melting and freezing, interaction of sea-ice floes, and also evaluating the velocity of two-dimensional ice-drift;

3: Sea-ice advection by this drift velocity;

4: Advection of 3D scalars, such as temperature and salinity;

5: Vertical diffusion of the scalars with sources due to heating, ice melting/freezing, and so on;

6: Dynamics of 3D horizontal current velocity;

7: Solving the SLAE for the sea level;

8: Evaluating the vertical velocity.

---

The matrix of the System of Linear Algebraic Equations (SLAE) is sparse and it contains 19 non-zero diagonals that correspond to adjacent mesh nodes within a finite element. The matrix does not vary in time and it is precomputed before the time step loop. The most time-consuming steps for the sequential code version were: 3D advection of scalars, 2D advection of sea-ice fields and solving the SLAE for the sea level. The simple Characteristic-Galerkin Scheme (Zienkiewicz and Taylor, 2000)
80  is used for the 3D and 2D advection terms.
Sea ice is considered to be an ensemble of multiple floes with some thickness distribution. This distribution is approximated by the discrete one with 15 fixed thickness values The local 1D sea ice thermodynamics is based on the 0-layer model (Semtner, 1976; Parkinson and Washington, 1979) with some modifications in lateral melting and surface albedo (Yakovlev, 2009). There are 14 categories of ice thickness (gradations), including zero thickness (open water). the mechanical redistribution and
85  the ice strength are identical to the CICE (Hunke et al., 2013). The elastic-viscous plastic scheme (EVP; Danilov et al., 2015) with modification for the relaxation time scales (Wang et al., 2016) is used for the sea ice dynamics (see also the Appendix 3 in Koldunov et . Sea ice is described by distribution of its compactness (concentration) for each gradation and ice volume for each gradation(excluding water). In addition, snow-on-ice volume for each gradation is evaluated. Therefore, there are 43 2D sea-ice scalars: ice and snow volume for 14 gradations and sea-ice compactness for 15 ones (including water). Because there are 39
90  vertical layers in an ocean component, the set of all of the sea-ice data is comparable to a single 3D scalar.

The tested version of the model has a spatial resolution of $0.036°$E, $0.011°$N, which is between $1.0$ and $1.3$ km along parallels and $1.2$ km along a meridian. The number of 2D grid nodes is $500 \times 500$, and only 33% of them are "wet" (84542). The time step is 100 s. The vertical step is 5 m up to 150 m deep and then 10 m up to 300 m. In fact, in the bathymetry data (ETOPO Amante and Eakins (2009)) (ETOPO; Amante and Eakins, 2009) the deepest point of the sea more shallow than it
95  really is, which reduces the actual maximum depth to 240 m. Comparison of available bathymetry data for the White Sea is given in Chernov and Tolstikov (2020) in table 1.

4

## 4 Organization of the calculations in non-parallel code

Computations in the ocean and sea-ice components are performed using three-dimensional arrays, such as $a(i, j, k)$ or $b(i, j, m)$, where $i, j$ represent the horizontal grid indices, $k$ represents the depth-layer, and $m$ represents the ice gradation. The differential operators are local: only neighbouring grid nodes—that is, $a(i \pm 1, j \pm 1, k)$—are used.

Typical differential operators in the ocean component are organized as shown in algorithm 2, where $N_x = 500, N_y = 500$ and $K(i, j)$ is the number of depth layers. For land points, $K(i, j) = 0$, and $K(i, j) \in [3, 39]$ with approximate mean value 12 for the remaining "wet" points, see figure 1.

---

**Algorithm 2** Typical 3D calculation loop

---
1: **for** $j = 1, N_y$ **do**
2:    **for** $i = 1, N_x$ **do**
3:       **for** $k = 1, K(i, j)$ **do**
4:          $a(i, j, k) = ...$
5:       **end for**
6:    **end for**
7: **end for**

---

Differential operators in the ice component are shown in algorithm 3, where $M = 14$ or 15 is the number of ice gradations and $\mathrm{mask}(i, j)$ is the logical mask of wet points. The percentage of wet points is 33 %.

---

**Algorithm 3** Typical 2D calculation loop

---
1: **for** $j = 1, N_y$ **do**
2:    **for** $i = 1, N_x$ **do**
3:       **if** $\mathrm{mask}(i, j)$ **then**
4:          **for** $m = 1, M$ **do**
5:             $b(i, j, m) = ...$
6:          **end for**
7:       **end if**
8:    **end for**
9: **end for**

---

Note that arrays in Fortran are arranged in the column-major order, so the first index $i$ is linear in memory. The presented arrangement of indices is common for ocean models ~~, see for example NEMO Madec et al. (2015)~~(see, for example, NEMO documentation; N. The loops arrangement is utilized from the original code. Although another arrangement may be more efficient, it does not affect the parallelization approach given later on. In spite of the fact that inner loop does not have stride-1 access, we can speculate that it allows for possible automatic vectorization over $m$ index and corresponds to minimal control flow interruptions due to *false* $\mathrm{mask}(i, j)$ values.
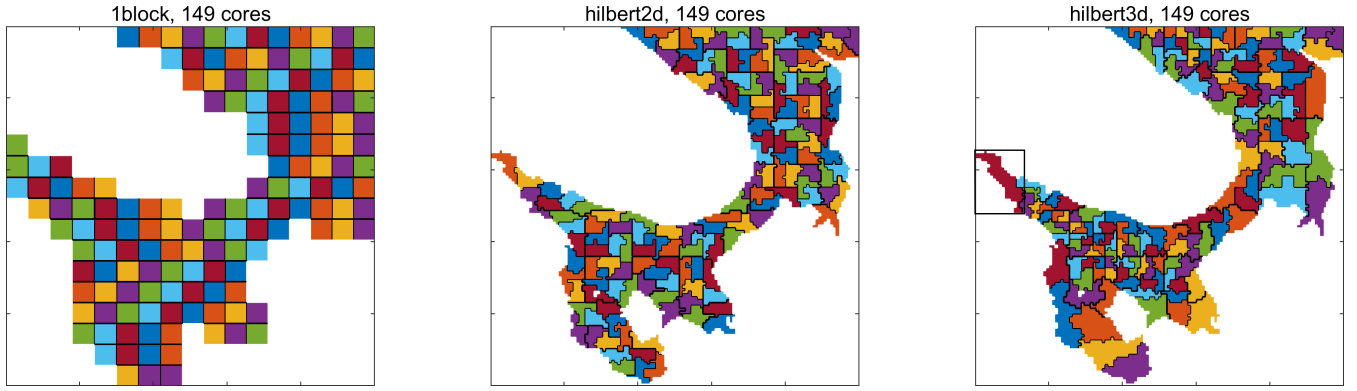
**Figure 2.** Three types of partition. Different processor cores are separated by a black line. Hilbert partitions are based on a grid of $n_b \times n_b = 128 \times 128$ blocks. Colours can repeat. Black rectangle in figure for hilbert3d partition corresponds to a "shared" array allocated for blocks belonging to a given CPU core.

## 5  Modifications of the non-parallel code

In this section we describe the partitioning algorithm of the model domain into subdomains, each corresponding to a CPU core, and subsequent modifications of the single-core calculations, which require only minor changes of algorithms 2 and 3. Grid

115  partition is performed in two steps: model domain is decomposed into small blocks and then these blocks are distributed over CPU cores in such a way that computational load imbalance is minimized. We utilize common grid partition for both sea-ice and ocean submodels, and provide theoretical estimates of the load imbalances resulting from the application of different weight functions in the balancing problem. Partition is calculated during the model initialization stage, as our balancing algorithm (Hilbert curves) is computationally unexpensive. Also, we guarantee that the partition is the same each time the model is run,

120  if parameters of the partitioner were not modified.

Computational domain $[1, N_x] \times [1, N_y]$ is separated into $n_b \times n_b$ blocks. If integer division is impossible, then block sizes are $n_x(i_b, j_b) = N_x \operatorname{div} n_b$, $n_y(i_b, j_b) = N_y \operatorname{div} n_b$, where $i_b, j_b \in [1, n_b]$ are horizontal indices of the blocks. The other points are distributed over the first blocks: $n_x(i_b, j_b) \mathrel{+}= 1$, where $i_b \in [1, N_x \bmod n_b]$, $j_b \in [1, n_b]$ and $n_y(i_b, j_b) \mathrel{+}= 1$, where $i_b \in [1, n_b]$, $j_b \in [1, N_y \bmod n_b]$. The set of indices corresponding to a block is denoted by $\Omega(i_b, j_b) = [i_s(i_b, j_b), i_e(i_b, j_b)] \times$

125  $[j_s(i_b, j_b), j_e(i_b, j_b)]$.

To formulate a balancing problem, we must assign weights of computational work to each block and then distribute them among $N_p$ available CPU cores in such a way that all cores have the same amount of work to do, or as close to this as possible, but provided that the "quality" of the partition is kept. Connectivity of subdomains (by subdomain we refer to a set of blocks belonging to a CPU core) or minimum length of the boundary can be chosen as possible criteria for the quality of a partition.

130  The weight for a block is the sum of weights corresponding to grid points in the range $\Omega(i_b, j_b)$. The following weights are

chosen for 2D and 3D computations, respectively:

$$w_{2d}(i,j) = \text{mask}(i,j), \tag{1}$$

$$w_{3d}(i,j) = K(i,j)/\text{mean}(K), \tag{2}$$

where "mean" operation is applied over wet points.

## 5.1 Trivial 1block partition

For a fixed $n_b$, one can find the number of "wet" blocks (i.e., blocks with at least one not-land point). In this partition, the number of cores $N_p$ is equal to the number of wet blocks and each CPU core gets exactly one block, see figure 2. Varying $n_b$, possible values of $N_p$ can be found.

## 5.2 Hilbert curve partition

For $n_b$ being a power of 2, the Hilbert curve connecting all the blocks can be constructed ~~Bader (2012)~~(Bader, 2012). This gives a one-dimensional set of weights that is balanced using the simplest algorithm. The sum of the blocks' weights on $p$ core is denoted by $W_p$. In spite of the fact that the Hilbert curve possesses the locality property (i.e., close indices on the curve correspond to close indices on the grid), it may not provide a partition into connected subdomains if there are a lot of land blocks. To overcome the problem of possible loss of connectivity, we perform the following optimization procedure, see algorithm 4.

---
**Algorithm 4** Optimization of partition
---
1: remove_not_connected_subdomains();
2: **for** $iter = 1, N_{iter}$ **do**
3:     balance_all_ranks();
4:     remove_not_connected_subdomains();
5: **end for**
---

Function `remove_not_connected_subdomains()` finds the connected subdomain with the maximum work for each CPU core and sends other blocks to neighbouring cores. Function `balance_all_ranks()` tries to send bordering blocks for each core to neighbouring cores to minimize the maximum work on both cores: $\max(W_p, W_{p'}) \to \min$, where $W_p, W_{p'}$ are for the work on the current CPU core and on a neighbouring core, respectively. The number of iterations is user-defined and we choose $N_{iter} = 15$, which is usually enough to reach convergence. Note that optimization does not guarantee to find a global optimum ~~.~~and function `remove_not_connected_subdomains()` may increase LI. Thus, we choose the iteration with the best balancing. The need for partitioning into connected subdomains comes from the intention to increase percentage of the wet points on CPU cores due to the data structure used; see the following section for a definition of the "shared" array.

The described algorithm performs partitioning into connected subdomains with Load Imbalance, which is

$$LI = 100\% \cdot \frac{\max(W_p) - \mathrm{mean}(W_p)}{\mathrm{mean}(W_p)},\ p \in [1, N_p], \tag{3}$$

not more than 10% in most cases. This is an acceptable accuracy because partitioning itself is not the main objective of the article.

~~Let us introduce~~ We have implemented two baseline partitions: hilbert2d (with weights $w_{2d}$) and hilbert3d (with weights $w_{3d}$), see figure 2. As one can see, hilbert2d divides the computational domain on quasi-uniform subdomains, while hilbert3d locates many CPU cores in high-depth regions and few cores in shallow water. Minimum and maximum number of blocks on a core can be found in tables 2, 3. When one of these partitions is applied to the whole coupled ocean-ice model, it balances one submodel and unbalances another. Table 2 shows that balancing of 2D computations ("LI 2D" $\rightarrow \min$) leads to imbalance in 3D computations ("LI 3D" $\approx 200\%$) and table 3 shows the opposite behaviour with "LI 2D" $\approx 300\%$. These values are close to the estimates given in appendix A and defined by the ratio between minimum, maximum and mean integer depth. The presented LI values imply a slowdown of one of the submodels by three to four times because LI increases runtime ($T$) compared to optimal one ($T_{opt}$) in the following way:

$$T = (LI + 1)T_{opt}. \tag{4}$$

A compromise for both submodels can be found by considering a combination of weights:

$$w_{2d3d} = w_{2d} + \gamma_0 w_{3d}, \tag{5}$$

where $\gamma_0 \approx 3$ is a ratio of run times for ocean and ice submodels on one CPU core. A partition of this type is denoted by hilbert2d3d. While this weight is optimal for "overlapping" computations of two code sections with different complexity, it is also the optimal weight for "non-overlapping" code sections (i.e., separated by blocking MPI exchanges). We show this in appendix B with corresponding estimates of LI for 2D (130%) and 3D (34%) computations.

## 5.3 Data structure and MPI exchanges

After partitioning has been performed, we get a set of blocks for each CPU core $p$, $I_p = \{(i_b, j_b)\}$. "Shared" data arrays are allocated for all blocks belonging to a CPU core with the following range of indices (excluding halo):

$$i_s^p = \min(i_s(I_p)), i_e^p = \max(i_e(I_p)), \tag{6}$$

$$j_s^p = \min(j_s(I_p)), j_e^p = \max(j_e(I_p)), \tag{7}$$

$$a(i_s^p : i_e^p, j_s^p : j_e^p, :). \tag{8}$$

~~The~~ An example of the shared array size is shown by rectangle in figure 2 for a particular CPU core. We introduce a mask of grid points belonging to a CPU core ($\mathrm{mask}_p(i,j)$). Correspondence between blocks, shared array and the mask is clarified in figure 3. Introducing this mask does not increase the complexity of the algorithms because the mask of the wet points already
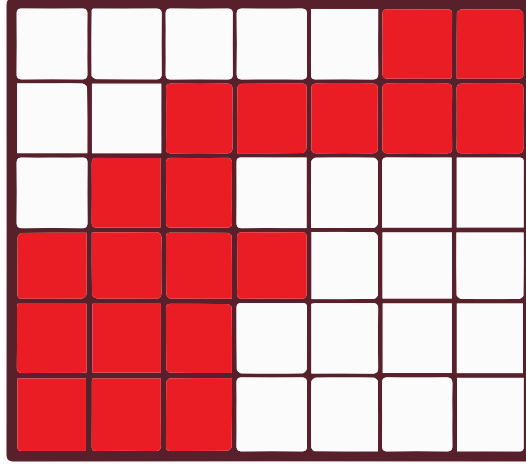
**Figure 3.** Blocks belonging to a CPU core in colour and borders of allocated array by thick line. $\mathrm{mask}_p(i,j) = 1$ in coloured blocks and $0$ elsewhere.

exists in the original code and it is simply modified. Finally, only minor modifications of the original loops are necessary:

$$1, N_x \to i_s^p, i_e^p, \tag{9}$$

185
$$1, N_y \to j_s^p, j_e^p, \tag{10}$$

$$K(i,j) \to K(i,j) \cdot \mathrm{mask}_p(i,j), \tag{11}$$

$$\mathrm{mask}(i,j) \to \mathrm{mask}(i,j) \cdot \mathrm{mask}_p(i,j). \tag{12}$$

Usually, see ~~Dennis (2007, 2003)~~Dennis (2007, 2003), arrays are allocated for each block separately. This has the following advantages:

190
- More efficient cache usage;

- If the number of blocks is large enough to get proper balancing, then there is no need for $\mathrm{mask}_p(i,j)$, thus giving advance in vectorization and so on.

It also introduces some drawbacks:

- Overheads for copying block boundaries (small blocks like $4 \times 4$ are prohibited);

195
- Many modifications of the original code are necessary, especially in service routines, I/O, and so on.

Consequently, the main strength of our approach is the ability to incorporate balancing while keeping the original program code as simple as for the trivial 1block partition. We expect that a "shared" array may be not optimal for near-land CPU cores with $\sim 20\%$ of wet points because of non-efficient cache usage. An example of such a core is shown by the rectangle in figure 2. An experimental study of runtime dependence on % of wet points will be carried out.

**Table 1.** The model with 1block partition; "$n_b \times n_b$" is the grid of blocks; "LI 2D" and "LI 3D" are load imbalances (3) for weights $w_{2d}$ and $w_{3d}$, respectively. "LI iceadvect" and "LI advect3D" are LIs computed based on the runtime of corresponding functions without exchanges; "days / 24 hours" is the number of computed days for one astronomical day.

| CPU cores | 1 | 32 | 78 | 149 | 306 | 595 | 993 |
|---|---|---|---|---|---|---|---|
| $n_b \times n_b$ | $1^2$ | $7^2$ | $12^2$ | $17^2$ | $26^2$ | $38^2$ | $50^2$ |
| LI 2D, % | 0 | 93 | 62 | 53 | 37 | 28 | 17 |
| LI iceadvect, % | 0 | 80 | 57 | 50 | 40 | 30 | 19 |
| LI 3D, % | 0 | 341 | 317 | 380 | 313 | 278 | 274 |
| LI advect3D, % | 0 | 339 | 324 | 340 | 349 | 290 | 291 |
| days / 24 hours | 8 | 79 | 219 | 360 | 864 | 1763 | 2556 |

200  Borders of blocks neighbouring with other CPU cores are sent using MPI. The following optimizations are applied to reduce the exchange time:

– All blocks' boundaries adjacent to a given CPU core are copied to a single buffer array, which is sent in one `MPI_Send` call.

– If possible, a diagonal halo exchange is included into cross exchanges with extra width.

205  – There is an option to send borders of two or more model fields in one `MPI_Send` call. ~~These three bullets reduce latency cost in many cores.~~

– Borders in the sea component are sent up to $K(i,j)$ depth (i.e., only the actually used layers are transmitted). ~~This~~

The first three bullets reduce latency cost in many cores and the final bullet reduces bandwidth limitations.

### 5.4  Parallel solver of the SLAE

210  As we have already mentioned, the time-implicit equation for the free surface is reduced to a SLAE with sparse 19-diagonal matrix. This is solved by a parallel implementation of Bicgstab algorithm preconditioned by block-ILU(0) with overlapping blocks, see ~~Saad (2003)~~Saad (2003). ILU(0) preconditioner preserves the 19-diagonal matrix structure, where matrix blocks are defined for each CPU core and correspond to wet points plus a band of border points of width 2. Because blocks are defined by the partition, the convergence rate depends on the number of CPU cores. Nevertheless, we have found that in the range from

215  1 to 996 CPU cores, it is sufficient to perform 6 to 10 iterations in order to reach the relative residual $\|Ax - b\| / \|b\| \leq 10^{-6}$.

**Table 2.** Same as table 1, but for hilbert2d partition; min and max operations are applied over CPU cores; column "estimate" shows theoretical LI given in appendix A.

| CPU cores | 1 | 32 | 78 | 149 | 306 | 595 | 993 | estimate |
|---|---|---|---|---|---|---|---|---|
| $n_b \times n_b$ | $1^2$ | $64^2$ | $128^2$ | $128^2$ | $128^2$ | $128^2$ | $128^2$ | |
| min blocks | 1 | 41 | 62 | 34 | 17 | 8 | 5 | |
| max blocks | 1 | 61 | 103 | 54 | 30 | 17 | 11 | |
| min % of wet points | 33 | 20 | 23 | 28 | 27 | 22 | 19 | |
| LI 2D, % | 0 | 9 | 8 | 7 | 4 | 12 | 28 | 0 |
| LI iceadvect, % | 0 | 11 | 19 | 12 | 10 | 19 | 27 | |
| LI 3D, % | 0 | 147 | 195 | 205 | 213 | 222 | 255 | 225 |
| LI advect3D, % | 0 | 145 | 200 | 217 | 242 | 235 | 264 | |
| days / 24 hours | 8 | 129 | 278 | 523 | 890 | 1826 | 2511 | |

**Table 3.** Same as table 2, but for hilbert3d partition.

| CPU cores | 1 | 32 | 78 | 149 | 306 | 595 | 993 | estimate |
|---|---|---|---|---|---|---|---|---|
| $n_b \times n_b$ | $1^2$ | $64^2$ | $64^2$ | $128^2$ | $128^2$ | $128^2$ | $128^2$ | |
| min blocks | 1 | 14 | 5 | 11 | 5 | 2 | 1 | |
| max blocks | 1 | 186 | 80 | 155 | 93 | 55 | 39 | |
| min % of wet points | 33 | 25 | 22 | 22 | 22 | 26 | ~~21~~ 25 | |
| LI 2D, % | 0 | 237 | 288 | 268 | 312 | 311 | ~~305~~ 313 | 300 |
| LI iceadvect, % | 0 | 238 | 297 | 298 | 373 | 337 | ~~329~~ 328 | |
| LI 3D, % | 0 | 5 | 7 | 3 | 12 | 15 | ~~46~~ 29 | 0 |
| LI advect3D, % | 0 | 31 | 21 | 18 | 16 | 23 | ~~48~~ 36 | |
| days / 24 hours | 8 | 131 | 338 | 691 | 1216 | 2232 | ~~3130~~ 3143 | |

## 6 Numerical experiments

Our experiments were performed on the cluster of Joint Supercomputer Center of the Russian Academy of Sciences[1]. Each node includes two 16-core processors Intel Xeon E5-2697Av4 (Broadwell). The software code was compiled by the Intel Fortran Compiler ifort 14.0.1 with the optimization option -O2. ~~Simulations~~ Low-core simulations were performed for three model days (2592 time steps). The model on 993 CPU cores is launched for 30 days, with subsequent rescaling of the results. During the first day, we call an `MPI_Barrier` function to measure performance of particular procedures with and without exchanges. During the last two days, an `MPI_Barrier` is omitted and overall performance is assessed. The ~~model is launched~~

---

[1] http://www.jscc.ru/

**11**

**Table 4.** Same as table 2, but for hilbert2d3d partition; "estimate" is given in appendix B.

| CPU cores | 1 | 32 | 78 | 149 | 306 | 595 | 993 | estimate |
|---|---|---|---|---|---|---|---|---|
| $n_b \times n_b$ | $1^2$ | $64^2$ | $64^2$ | $128^2$ | $128^2$ | $128^2$ | $128^2$ | |
| min blocks | 1 | 16 | 6 | 14 | 6 | 3 | 2 | |
| max blocks | 1 | 112 | 53 | 104 | 64 | 38 | 23 | |
| min % of wet points | 33 | 22 | 14 | 23 | 23 | 27 | 24 | |
| LI 2D, % | 0 | 95 | 126 | 131 | 130 | 142 | 139 | 130 |
| LI iceadvect, % | 0 | 119 | 138 | 150 | 156 | 150 | 145 | |
| LI 3D, % | 0 | 19 | 27 | 26 | 27 | 41 | 66 | 34 |
| LI advect3D, % | 0 | 24 | 34 | 32 | 29 | 48 | 72 | |
| days / 24 hours | 8 | 180 | 403 | 811 | 1615 | 2718 | 3463 | |



**Figure 4.** Speedup compared to one core for two functions: `advect3D` and `iceadvect`. Solid/dashed lines correspond to measurements with/without ~~MPI-exchanges~~MPI exchanges, respectively. Different partitions are shown in colour (1block, hilbert2d, and hilbert3d).

~~on 993 CPU cores for 30 days, with subsequent rescaling of the results. The~~ number of cores for tests are guided by the 1block partition method, which is highly restricted in the allowable number of cores. We first show how the most time-consuming functions corresponding to ocean and ice submodels accelerate for three partitions: 1block, hilbert2d and hilbert3d (see figure 2). We then study overall performance of the model using four partitions, including hilbert2d3d with combined weights (5).

The maximum grid size of blocks for our model is $n_b \times n_b = 128 \times 128$ because the MPI exchange width is limited by the block size, while the SLAE solver requires exchange of width 2. Note that due to the data structure that we used, the
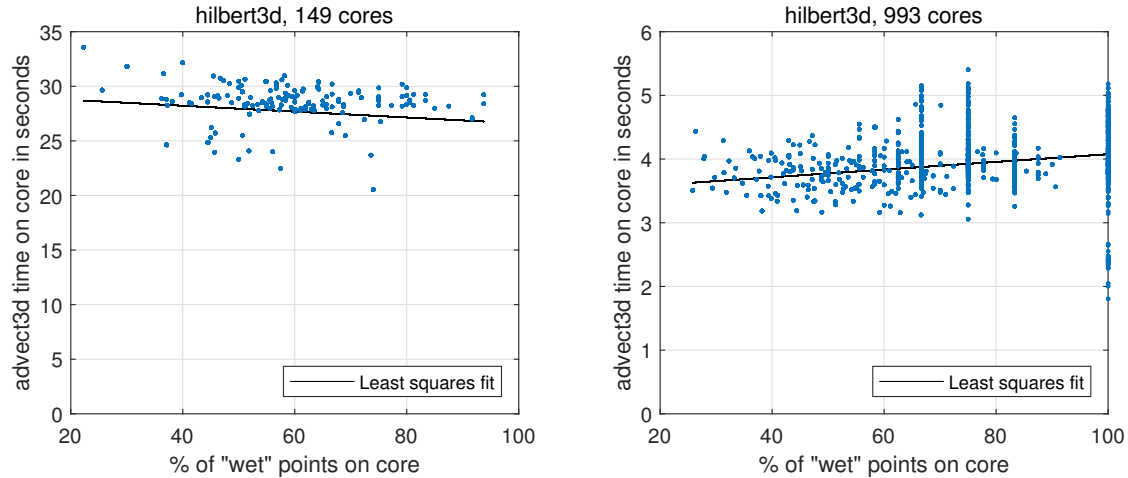
**Figure 5.** Scatter plot: percentage of wet points on a core – `advect3D` runtime (without MPI exchanges, for 1 model day). Each point corresponds to one CPU core. Figures correspond to hilbert3d partition with different numbers of CPU cores.

performance of Hilbert-type partitions is almost insensitive to $n_b$ at moderate number of CPU cores. Nevertheless, $n_b$ may
230 be tuned by hand to decrease the complexity of the partition optimization procedure or to increase the percentage of the wet points on a core. In runs with many CPU cores, we use the maximum available number of blocks to get better balancing. The parameters that we used in the experiments are given in tables 1, 2, 3, 4.

## 6.1 Speedup of scalar and ice advection

Advection of scalars (`advect3D`, depth-dependent) and ice (`iceadvect`, depth-independent) are the most time-consuming
235 procedures in ocean and ice submodels, respectively. In the following, we will show that hilbert3d partition is appropriate for `advect3D` and hilbert2d for `iceadvect`.

Speedup for mentioned procedures is given in figure 4. Dashed lines correspond to measurements of code sections between MPI exchanges and show how pure computations accelerate. Pure computations in `advect3D` accelerate linearly for hilbert3d partition, while pure computations in `iceadvect`—superlinearly for the hilbert2d partition. ~~We explain superlinearity by~~
240 ~~better cache usage.~~ Superlinearity (parallel efficiency is greater than 100% when "doubling" the number of cores) occurs in the range of cores $1 - 149$, and we have investigated the possible reasons for this using the open access Intel Advisor profiler. We have found that on 1 core the most time demanding memory requests are the DRAM data upload, while on 149 cores most memory requests correspond to the L1 cache. Thus, superlinearity can be partially explained by the better cache utilization when the number of cores increases. Note that more local memory access pattern (MAP) can decrease the limitations caused
245 by the memory requests and can be achieved by incorporating stride-1 access for the inner loop indices, but we leave this point for optimizations in the future.
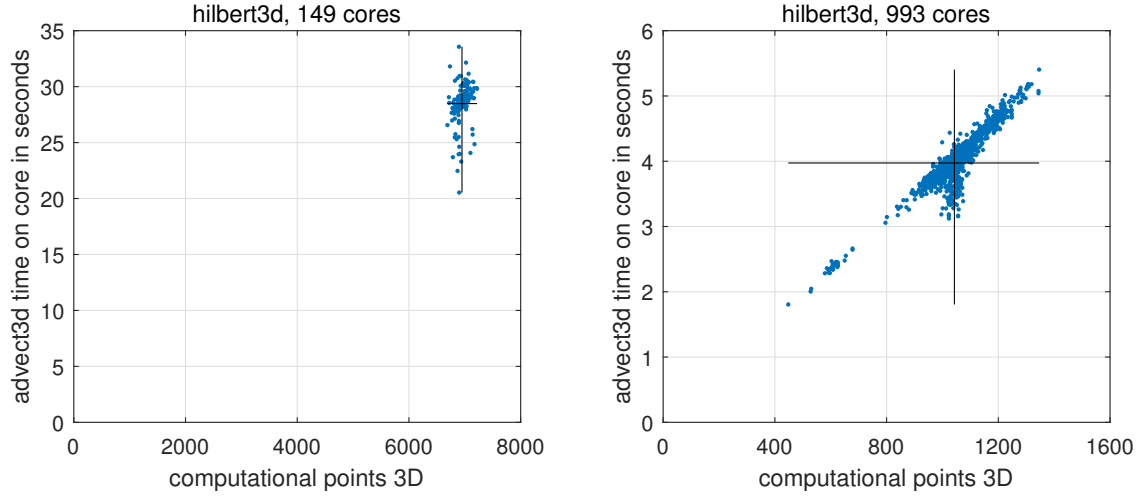
**Figure 6.** Scatter plot: number of ~~blocks on a core – number of~~ computational points for 3D calculations <u>on a core – `advect3D` runtime</u> <u>(without MPI exchanges, for 1 model day)</u>. Each point corresponds to one CPU core. <u>Solid lines show the average values along x and y axes.</u> Figures correspond to hilbert3d partition with different numbers of CPU cores.

<u>Speedup including MPI exchanges is shown in figure 4 with solid lines.</u> The function `advect3D` is only slightly limited by MPI exchanges on 993 cores: its speedup on the partition hilbert3d falls from ~~865 to 615~~ <u>940 to 640</u> when exchanges are accounted for. Meanwhile, `iceadvect` loses speedup from 1540 to 576 after accounting for exchanges on hilbert2d partition.

250 Both functions have identical number of exchanges, but `advect3D` is more computationally expensive. Consequently, we explain worse performance of `iceadvect` by lower ratio of number of operations to the number of points to exchange. Similar bottleneck due to exchanges in 2D dynamics is reported in ~~Koldunov et al. (2019a)~~<u>Koldunov et al. (2019a)</u>. The hilbert2d partition has a slight advantage (about 15–20%) over the 1block partition for both functions (see solid lines). In total, as we expected, the hilbert3d partition is suitable for `advect3D` function, and its acceleration is two to three times more efficient

255 than when 1block/hilbert2d partitions are used. Also, hilbert2d/1block partitions show two to four times faster `iceadvect` function compared to hilbert3d partition. The different accelerations are strongly connected to balancing of computations. To check partition-based ("LI 3D" and "LI 2D") and runtime-based ("LI advect3D" and "LI iceadvect", correspondingly) Load Imbalance for 3D and 2D computations, see tables 1, 2, 3. Note that theoretical and practical LI are moderately close to each other, which confirms our choice of weights (1), (2) for these functions. Also note that the data structure and organization of

260 the calculations are appropriate for load balancing.

Further analysis reveals that the runtime-based LI could be 4–25% more than the partition-based one, see tables 2, 3. This may be connected to overheads introduced by non-efficient organization of memory. We allocate a shared array for all blocks belonging to a CPU core and near-land cores may have only 20 % of the wet points (see tables 2, 3), which can lead to an increase in cache misses. Figure 5 shows a scatter plot for % of the wet points vs `advect3D` runtime without exchanges
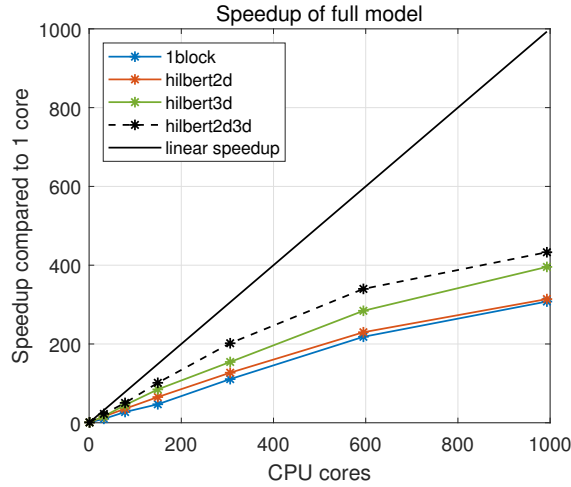
**Figure 7.** Speedup compared to one core for the full model. Different partitions shown in colour (1block, hilbert2d, hilbert3d, hilbert2d3d).

265  (each point corresponds to some CPU core). One can clearly see that on a moderate number of cores (149) the computations are limited by the core with the smallest % of the wet points, which has the maximal runtime. ~~However, there is no drastic dependence of runtime on % of wet points. This may mean that the number of operations per one array element in this model is large enough, thus the data structure plays a moderate role. In particular, the data structure~~ Figure 6 additionally shows that spread in runtime cannot be explained by the difference in the number of computational points, i.e. partitioning algorithm

270  works well for 149 cores. Although organization of the calculations may slightly limit model efficiency on a moderate number of cores, it does not limit the model efficiency on 993 cores~~(see figure 5), where~~, where major part of the `advect3D` ~~runtime suffers from imperfect balancing: cores with small number of blocks usually fall into 100% of wet points and has a wide range of run times, approximately from 1 to 5 seconds. Figure ?? additionally shows that on 993 cores the balancing is limited for processors, where the number of blocks per core is small. Proper balancing of 3D computations by 2D partitioning implies that~~

275  ~~the number of blocks per core should be in a wide range, while the minimum number should not be close to 1 (see~~ runtime spread is explained by the ~~left-hand panel in figure ??~~imperfect balancing (see figure 6), but not the data structure (see figure 5). Stagnation of the balancing procedure is evident from the fact that ~~we have only 84542 surface wet points, which corresponds to patches of size~~ the minimum number of blocks located on a CPU core is 1 for 993 cores, see table 3. Note that computational subdomain corresponding to one CPU core is small enough: on average, it has $9 \times 9$ ~~(~~horizontal points with 12 vertical levels ~~)~~

280  ~~on~~ for 993 cores~~, on average~~.

## 6.2 Speedup of the full model

The coupled ocean-ice model is launched on the same CPU cores for both submodels with the common horizontal partition. Input/output functions are sequential and utilize gather-scatter operations, which are given by our library of parallel exchanges. Speedup for the full model compared to one CPU core is given in figure 7. Maximum speedup, approximately 430, corresponds
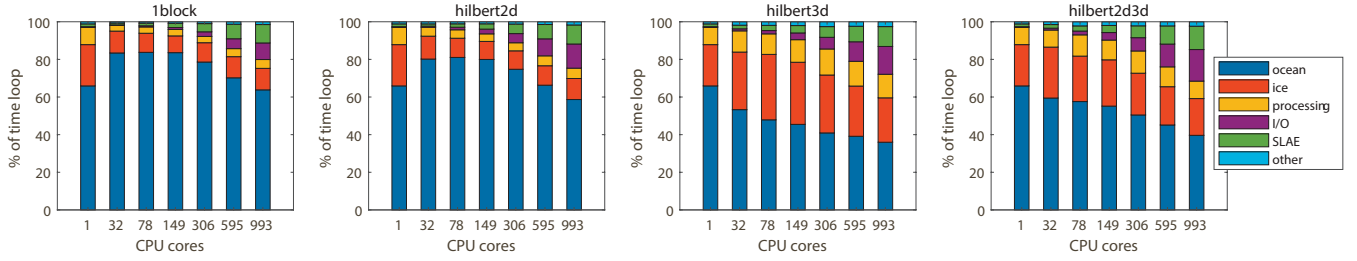
**Figure 8.** Relative contribution of different code sections to runtime; "ocean" – all procedures corresponding to ocean submodel including `advect3D`, "ice" – ice submodel including `iceadvect`, "processing" – computation of statistics, "I/O" – input/output with scatter-gather functions; "SLAE" – matrix inverse and RHS preparation; "other" – simple service procedures.
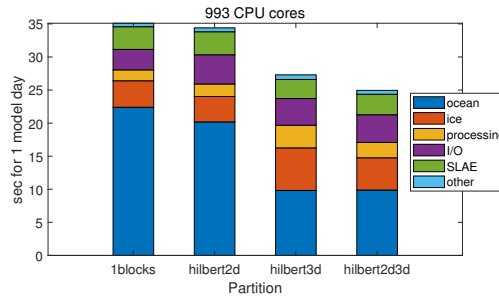


**Figure 9.** Absolute contribution of different code sections to runtime on 993 CPU cores.

to the partition with combined weights (hilbert2d3d) on 993 cores. Compared to the simplest partition (1block), hilbert2d3d model is 115% faster on 149 cores and 40% faster on 993 cores. Partition hilbert2d3d also gives an advantage over partitions balancing purely 2D and 3D computations (hilbert2d, hilbert3d). On 993 CPU cores, the parallel exchanges in this model have the following contribution to the runtime: 20% for boundary exchanges, 18% for gather-scatter and 6.5% for Allreduce.

The relative contribution of different code sections to runtime is given in figure 8. In the case of perfect scaling of all procedures, the relative contribution must be the same as the number of CPU cores rises. For partitions 1block and hilbert2d, we see a slowdown of the ocean component. Partition hilbert3d suffers from the slowdown of the ice component. Finally, the closest preservation of time distribution is found for hilbert2d3d model. We did not pay much attention to code section "processing" because, although it accelerates, its computational cost could be reduced. Section "I/O" gradually saturates due to gather-scatter operations, which consume 85% of I/O runtime on 993 CPU cores. The new parallel solver ("SLAE") has fast convergence and low computational cost, but suffers from Allreduce operations: in our implementation each iteration demands five `MPI_Allreduce` calls, which account for 60% of "SLAE" code section runtime on 993 cores.

Absolute values of code sections' runtime for 1 model day are shown in figure 9. In comparison to hilbert3d partition, combining of weights (hilbert2d3d) reduces the cost of the ice component, while keeping ocean component almost without changes (see also tables 4, 3 for load imbalance values). In addition, section "processing" reduces its runtime because it contains

300    many not fully optimized service functions that are sensitive to stretching of the horizontal area covered by a CPU core: such stretching is done by the hilbert3d partitioner.

Simulated years per wall-clock day (SYPD) for the best configuration (hilbert2d3d, 993 cores) is $3463/365 \approx 9.5$, see table 4. A direct comparison with other coupled ocean-ice models cannot be achieved because our configuration is rare. However, we can rescale the performance (rSYPD) of time step efficiency of the global models in the following way:

$$305 \quad rSYPD = SYPD \frac{N_{mesh}}{N_{mesh}^{FEMAO}} \frac{\Delta t^{FEMAO}}{\Delta t} \frac{N_p^{FEMAO}}{N_p}, \tag{13}$$

where we take into consideration different numbers of horizontal mesh wet points ($N_{mesh}$), CPU cores ($N_p$) and time step ($\Delta t$), but we neglect different numbers of vertical levels and differences in formulation of the ice dynamics. As follows from table 5, rSYPD is of the order of 10 for all of the ocean-ice models that we have presented, including FEMAO. While this characteristic cannot rate models over their efficiency, we argue that our parallel configuration is comparable to existing parallel ocean-ice
310    models.

**Table 5.** Efficiency of time step loop for FEMAO model compared to global ocean-ice models. Rescaled SYPD (rSYPD, (13)) accounts for difference in the number of horizontal mesh points, CPU cores and time step. Original values are published in ~~Koldunov et al. (2019a); Huang et al. (2016); Ward (2016)~~Koldunov et al. (2019a); Huang et al. (2016); Ward (2016), but we took our values directly from table 3 in ~~Koldunov et al. (2019a)~~Koldunov et al. (2019a).

| Model | Mesh points $\cdot 10^6$ | Cores | Time step, s | SYPD | rSYPD |
|---|---|---|---|---|---|
| POP | 5.8 | 16875 | 173 | 10.5 | 24.4 |
| FESOM2/STORM | 5.6 | 13828 | 600 | 15.9 | 12.5 |
| NEMO | 0.9 | 3840 | 1440 | 25.3 | 4.8 |
| MOM5.1 | 0.9 | 3840 | 1800 | 21.6 | 3.3 |
| FESOM2/farc | 0.6 | 2304 | 900 | 56.2 | 19 |
| FEMAO | 0.085 | 993 | 100 | 9.5 | 9.5 |

## 7    Conclusions

In this paper, we present a relatively simple approach to accelerate the FEMAO ocean-ice model based on rectangular structured grid with advances of load balancing. The modifications that had to be introduced into the program code are identical to those that were required by the simplest decomposition on squares. The only demand on the model to be accelerated by this
315    technique is marking computational points by a logical mask. In the first step, we utilize the common partition for ocean and ice submodels. For a relatively "small" model configuration, $500 \times 500$ horizontal points, we reach parallel efficiency of 60 % for particular functions (3D scalar advection using 3D-balancing approach and 2D ice advection using 2D-balancing approach) and 43% for the full model (using combined weight approach) on 993 CPU cores. We show that balancing the 3D computations

leads to unbalanced 2D computations, and vice versa. Consequently, further acceleration may be achieved by performing

320 computations of 2D and 3D components on distinct groups of CPU cores with different partitions. Nevertheless, high parallel efficiency of 3D scalar advection itself is a great advance for future applications of the model, especially for the version with a pelagic ecology submodel ~~Chernov et al. (2018)~~(Chernov et al., 2018), where more than 50 3D scalars (biogeochemical concentrations) are added to the thermohaline fields.

Note that while the parallel approach that we have presented here can be implemented into the model in relatively simple

325 way, the code of the library of parallel exchanges can be rather complex (see Supplements).

## Appendix A: Estimating the load imbalance for hilbert2d and hilbert3d partitions

330 Let us introduce two functions of bathymetry (defined by integer depth $K(i,j)$) with the corresponding values for our model:

$$\rho_{max}(K) = \frac{\max(K)}{\text{mean}(K)} = \frac{39}{12} = 3.25, \tag{A1}$$

$$\rho_{min}(K) = \frac{\text{mean}(K)}{\min(K)} = \frac{12}{3} = 4, \tag{A2}$$

here and below, "mean", "min" and "max" operations correspond only to wet points. These values define how balancing of 2D computations affects 3D computations imbalance, and vice versa. Let $S$ and $V$ be sets of surface and ocean points,

335 correspondingly; $S_p$ and $V_p$ be sets of these points belonging to a CPU core $p$; $|\cdot|$ be the number of points in a set. The number of 3D points can be expressed via 2D ones: $|V| = \sum_{\{i,j\} \in S} K(i,j) = |S| \cdot \text{mean}_{\{i,j\} \in S} K(i,j)$.

When balancing of 2D computations is used (hilbert2d), surface points are distributed among processors in roughly equal size ($|S_p| = |S|/N_p$). Then, for 3D computations, the ratio of maximum work to mean work among cores is defined as:

$$\frac{W_{max}}{W_{mean}} = \frac{\max_p(|V_p|)}{\text{mean}_p(|V_p|)} = \frac{\max_p(\text{mean}_{\{i,j\} \in S_p} K(i,j))}{\text{mean}_p(\text{mean}_{\{i,j\} \in S_p} K(i,j))} \approx \rho_{max}, \tag{A3}$$

340 and the corresponding load imbalance is

$$LI = \frac{W_{max} - W_{mean}}{W_{mean}} = \rho_{max} - 1 = 225\%. \tag{A4}$$

When balancing of 3D computations is used (hilbert3d), ocean points are distributed among processors in roughly equal size ($|V_p| = |V|/N_p$). Then, for 2D computations, the ratio of maximum work to mean work is defined as:

$$\frac{W_{max}}{W_{mean}} = \frac{\max_p(|S_p|)}{\text{mean}_p(|S_p|)} = \frac{\max_p(|V_p|/\text{mean}_{\{i,j\} \in S_p} K(i,j))}{\text{mean}_p(|S_p|)} = \frac{|V|/|S|}{\min_p(\text{mean}_{\{i,j\} \in S_p} K(i,j))} \approx \rho_{min}, \tag{A5}$$

18

345 and the corresponding load imbalance is

$$LI = \frac{W_{max} - W_{mean}}{W_{mean}} = \rho_{min} - 1 = 300\%. \tag{A6}$$

## Appendix B: Finding the optimal weight for non-overlapping 2D and 3D calculations

Let $W$ be a full computational work and let it be distributed between 3D ($W^{3d}$) and 2D ($W^{2d}$) computations with ratio $\gamma_0$: $W = W^{2d} + W^{3d} \sim (1 + \gamma_0)W^{2d}$. Our goal is to find weight function $w(i,j)$, which corresponds to minimal joint (2D and 3D)

350 Load Imbalance. We use the notation presented in the previous appendix and we define the "number of computational points corresponding to weight": $|V^w| = \sum_{\{i,j\} \in S} w(i,j)$.

Assuming equipartition with respect to this weight ($|V_p^w| = |V^w|/N_p$), we can derive LI for 2D calculations:

$$\frac{W^{2d}_{max}}{W^{2d}_{mean}} = \frac{\max_p |S_p|}{\operatorname{mean}_p |S_p|} \approx \rho_{min}(w), \tag{B1}$$

and for 3D calculations:

355
$$\frac{W^{3d}_{max}}{W^{3d}_{mean}} = \frac{\max_p |V_p|}{\operatorname{mean}_p |V_p|} = \frac{\max_p(|S_p| \cdot \operatorname{mean}_{\{i,j\} \in S_p} K(i,j))}{\operatorname{mean}_p(|S_p| \cdot \operatorname{mean}_{\{i,j\} \in S_p} K(i,j))} = \frac{\max_p \left( \frac{\operatorname{mean}_{\{i,j\} \in S_p}(K(i,j))}{\operatorname{mean}_{\{i,j\} \in S_p}(w(i,j))} \right)}{\operatorname{mean}_p \left( \frac{\operatorname{mean}_{\{i,j\} \in S_p}(K(i,j))}{\operatorname{mean}_{\{i,j\} \in S_p}(w(i,j))} \right)} \approx \rho_{max}(K/w). \tag{B2}$$

Finally, assuming that 2D and 3D computations are non-overlapping (i.e., the maximum work is under summation), "Load Imbalance" for the full model:

$$LI(w) = \frac{W_{max} - W_{mean}}{W_{mean}} = \frac{\rho_{min}(w) + \gamma_0 \rho_{max}(K/w)}{1 + \gamma_0} - 1. \tag{B3}$$

For a given bathymetry $K(i,j)$, ratio $\gamma_0 = 3$ and special type of weight function $w(\gamma) = w_{2d} + \gamma w_{3d}$, $LI(w(\gamma))$ can be

360 plotted numerically for different values of $\gamma$, see figure A1. The minimum of this function corresponds to the choice $\gamma = \gamma_0 = 3$, and LI for 2D and 3D computations in this case are $130\%$ and $34\%$, respectively.
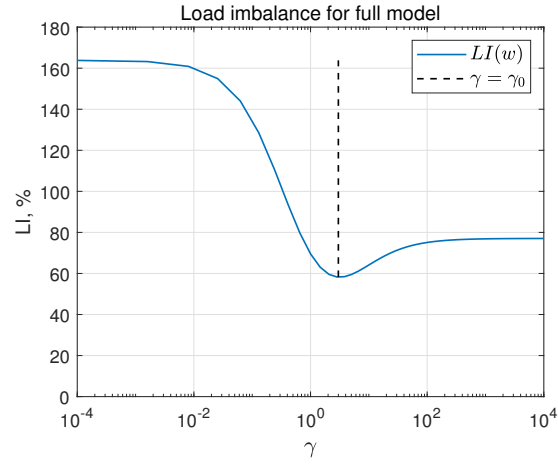
**Figure A1.** Load Imbalance for the full model $L(w)$, as a function of $w(\gamma) = w_{2d} + \gamma w_{3d}$; $\gamma_0 = 3$.

370

# References

Amante, C. and Eakins, B. W.: ETOPO1 arc-minute global relief model: procedures, data sources and analysis, Tech. rep., National Geophysical Data Center, 2009.

Bader, M.: Space-filling curves: an introduction with applications in scientific computing, vol. 9, Springer Science & Business Media, 2012.

Chaplygin, A. V., Dianskii, N. A., and Gusev, A. V.: Load balancing using Hilbert space-filling curves for parallel shallow water simulations, Vychislitel'nye Metody i Programmirovanie, 20, 75–87, 2019.

Chernov, I.: Numerical Modelling of large-scale Dynamics of the White Sea, Universal Journal of Geoscience, 1, 150–153, 2013.

Chernov, I. and Tolstikov, A.: The White Sea: Available Data and Numerical Models, Geosciences, 10, 463, https://doi.org/10.3390/geosciences10110463, 2020.

Chernov, I., Lazzari, P., and et. al.: Hydrodynamical and biogeochemical spatiotemporal variability in the White Sea: A modeling study, Journal of Marine Systems, 187, 23–35, 2018.

Danilov, S., Wang, Q., Timmermann, R., Iakovlev, N., Sidorenko, D., Kimmritz, M., Jung, T., and Schröter, J.: Finite-element sea ice model (FESIM), version 2, Geosci. Model Dev., 8, 1747–1761, 2015.

Dennis, J. M.: Partitioning with space-filling curves on the cubed-sphere, in: Proceedings International Parallel and Distributed Processing Symposium, pp. 6–pp, IEEE, 2003.

Dennis, J. M.: Inverse space-filling curve partitioning of a global ocean model, in: 2007 IEEE International Parallel and Distributed Processing Symposium, pp. 1–10, IEEE, 2007.

Filatov, N., Pozdnyakov, D., and et. al.: White Sea: its marine environment and ecosystem dynamics influenced by global change, Springer Science & Business Media, 2007.

Fox-Kemper, B., Alistair, and et. al.: Challenges and Prospects in Ocean Circulation Models, Frontiers in Marine Science, 6, 1–29, https://doi.org/10.3389/fmars.2019.00065, 2019.

Huang, X., Tang, Q., Tseng, Y., Hu, Y., Baker, A. H., Bryan, F. O., Dennis, J., Fu, H., and Yang, G.: P-CSI v1. 0, an accelerated barotropic solver for the high-resolution ocean model component in the Community Earth System Model v2. 0, Geoscientific Model Development, 9, 4209, 2016.

Hunke, E. C., Lipscomb, W. H., Turner, A. K., Jeffery, N., and Elliott, S.: CICE: the Los Alamos Sea Ice Model, Documentation and Software, Version 5.0, Los Alamos National Laboratory Tech. Rep. LA-CC-06-012, 2013.

Iakovlev, N.: On the calculation of large-scale ocean currents in the 'velocity-pressure' variables by the finite element method, Russian Journal of Numerical Analysis and Mathematical Modelling, 11, 383–392, 1996.

Iakovlev, N.: On the Simulation of Temperature and Salinity Fields in the Arctic Ocean, Izvestiya, Atmospheric and Oceanic Physics, 48, 86–101, https://doi.org/10.1134/S0001433812010136, 2012.

Karypis, G.: METIS, a software package for partitioning unstructured graphs, partitioning meshes, and computing fill-reducing orderings of sparse matrices version 4.0, http://glaros. dtc. umn. edu/gkhome/metis/metis/download, 1998.

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), Geoscientific Model Development, 12, 3991–4012, 2019a.

Koldunov, N. V., Danilov, S., Sidorenko, D., Hutter, N., Losch, M., Goessling, H., Rakowsky, N., Scholz, P., Sein, D., Wang, Q., et al.: Fast EVP Solutions in a High-Resolution Sea Ice Model, Journal of Advances in Modeling Earth Systems, 11, 1269–1284, 2019b.

Madec, G. et al.: NEMO ocean engine, Institut Pierre-Simon Laplace, 2015.

410   Parkinson, C. L. and Washington, W. M.: A large-scale numerical model of sea ice, Journal of Geophysical Research: Oceans, 84, 311–337, 1979.

Saad, Y.: Iterative methods for sparse linear systems, vol. 82, SIAM, 2003.

Semtner, A. J.: A model for the thermodynamic growth of sea ice in numerical investigations of climate, Journal of Physical Oceanography, 6, 379–389, 1976.

415   Wang, Q., Danilov, S., Jung, T., Kaleschke, L., and Wernecke, A.: Sea ice leads in the Arctic Ocean: Model assessment, interannual variability and trends, Geophysical Research Letters, 43, 7019–7027, 2016.

Ward, M.: Scalability of MOM5, NEMO and MOM6 on NCI's Raijin supercomputer, https://www.ecmwf.int/en/elibrary/16837-scalability-mom5-nemo-and-mom6-ncis-raijin-supercomputer, 2016.

Yakovlev, N.: Reproduction of the large-scale state of water and sea ice in the Arctic Ocean from 1948 to 2002: Part II. The state of ice and
420   snow cover, Izvestiya, Atmospheric and Oceanic Physics, 45, 478–494, 2009.

Zienkiewicz, O. and Taylor, R.: The Finite Element Method, 5th. Ed., Vol. 3: Fluid dynamics, Butterworth and Heinemann, Oxford, 2000.