

# Performance Analysis of the MPAS-Ocean Code using HPCToolkit and MIAMI

Gabriel Marin

February 11, 2014

MPAS-Ocean [4] is a component of the MPAS framework of climate models. MPAS-Ocean is an unstructured-mesh ocean model capable of using enhanced horizontal resolution in selected regions of the ocean domain. The code is publicly available for download [3] and comes with several input problems of different sizes corresponding to different simulation resolutions. In this initial study, we look at the per-core performance of version 2.0 of the MPAS-Ocean code. Our analysis was performed on a single node system with dual Intel Xeon E5-2690 CPUs, based on the Sandy Bridge micro-architecture. Each processor has 8 cores and a shared 20 MB L3 cache. We compiled the code with the Intel Fortran compiler 14.0.0 and optimization flags `-O3 -g`.

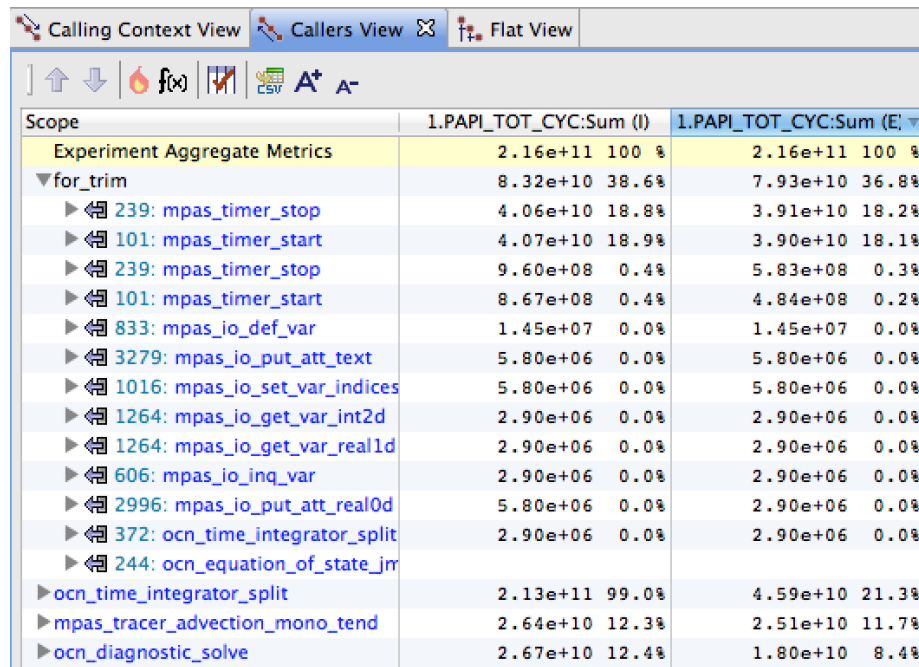


Figure 1: HPCToolkit time profile: the MPAS code spends 38.6% of time in the compiler intrinsic 'trim' invoked primarily from the software timer callipers.

Due to the small system size, we used an input problem of lower resolution (240km), and performed runs with one and four MPI processes. Figure 1 shows an initial HPCToolkit [1] time profile of the application running with a single MPI process. We use HPCToolkit's *callers* view and sort the results

scope	1.PAPI_TOT_CYC:Sum (I)	1.PAPI_TOT_CYC:Sum (E)
Experiment Aggregate Metrics	1.37e+11 100 %	1.37e+11 100 %
▶ocn_time_integrator_split	1.34e+11 98.3%	4.67e+10 34.2%
▶mpas_tracer_advection_mono_tend	2.66e+10 19.5%	2.51e+10 18.4%
▶ocn_diagnostic_solve	2.89e+10 21.2%	1.97e+10 14.4%
▶ocn_equation_of_state_jm_density	1.01e+10 7.4%	1.00e+10 7.3%
▶ocn_vel_coriolis_tend	5.27e+09 3.9%	5.27e+09 3.9%
▶ocn_fuperp	4.16e+09 3.0%	4.16e+09 3.0%
▶ocn_vmix_implicit	8.18e+09 6.0%	2.60e+09 1.9%
▶ocn_vert_transport_velocity_top	2.45e+09 1.8%	1.66e+09 1.2%
▶ocn_vel_forcing_windstress_tend	3.91e+09 2.9%	1.56e+09 1.1%

Figure 2: HPCToolkit time profile: the top time consuming routines in MPAS-Ocean after fixing the software timers.

by exclusive time, to identify the routines where the application spends the most time. Surprisingly, the application spends 38.6% of its execution time inside the Fortran compiler intrinsic `trim`. The *callers* view in HPCToolkit shows the places from where a routine is called, as well as a breakdown of the contribution of each call path to every metric. We notice that four call sites in two software timer calipers account for 38.5% of the entire time spent in routine `for_trim`.

Inspecting the timers code, we noticed repeated, unnecessary calls to the `trim` intrinsic for each invocation of the timer calipers. We changed the code to call the `trim` intrinsic once per caliper use, caching the result of its output. Figure 2 presents the updated HPCToolkit time profile after this code fix-up. The total execution time dropped by 36.5% from  $2.16e11$  to  $1.37e11$ . Routine `for_trim` accounts for 0.1% of the execution time in the new code version, thus, it does not show up in the list of top time consuming routines. We use the fixed version of the code as the baseline for further analysis and optimization.

## Performance analysis using MIAMI

Figure 3 shows a snapshot of several performance metrics computed by MIAMI [2]. MIAMI uses a modeling approach based on first-order principles to estimate the instruction schedule cost of an application on a particular architecture, to understand the data reuse patterns responsible for the highest number of data transfers between the various levels of the memory hierarchy, and to identify the memory accesses that cannot be effectively prefetched by streaming hardware prefetchers.

Scopes	CPU_Time	GainExtraRes	GainExtraIP	GainVectorize	RetiredUops	AddrGen	StackTmp
Experiment Aggregate Metrics	1.13e11 100.0	4.90e10 100.0	3.10e10 100.0	8.00e10 100.0	2.75e11 100.0	9.32e10 100.0	2.11e10 100.0
ocn_time_integration_split_mp_ocn_time_integrator_split	3.93e10 34.9%	1.20e10 24.4%	1.31e10 42.4%	3.28e10 41.0%	1.01e11 36.6%	6.70e10 71.9%	4.36e09 20.6%
mpas_tracer_advection_mono_mp_mpas_tracer_advection_m	1.80e10 16.0%	1.34e10 27.4%	2.78e09 9.0%	1.32e10 16.5%	5.65e10 20.5%	6.31e09 6.8%	2.38e09 11.3%
ocn_diagnostics_mp_ocn_diagnostic_solve	1.43e10 12.7%	8.89e09 18.1%	3.00e09 9.7%	9.42e09 11.8%	3.60e10 13.1%	4.88e09 5.2%	3.06e09 14.5%
ocn_equation_of_state_jm_mp_ocn_equation_of_state_jm_de	5.26e09 4.7%	4.33e09 8.8%	9.15e08 3.0%	1.68e09 2.1%	1.65e10 6.0%	1.08e09 1.2%	1.14e09 5.4%
ocn_diagnostics_mp_ocn_fuperp	4.12e09 3.7%	2.96e08 0.6%	1.53e09 4.9%	3.33e09 4.2%	1.13e10 4.1%	3.51e09 3.8%	1.46e09 6.9%
ocn_vel_coriolis_mp_ocn_vel_coriolis_tend	2.98e09 2.6%	1.11e09 2.3%	6.09e08 2.0%	2.17e09 2.7%	1.00e10 3.6%	2.81e09 3.0%	2.36e09 11.2%
ocn_vel_forcing_windstress_mp_ocn_vel_forcing_windstress	2.80e09 2.5%	7.69e07 0.2%	1.15e09 3.7%	1.98e09 2.5%	2.23e09 0.8%	7.59e08 0.8%	5.82e08 2.8%

Figure 3: MIAMI performance results: the top routines ranked by instruction schedule cost.

The metrics included in Figure 3 focus on the insight provided by MIAMI into the instruction schedule cost of the MPAS-Ocean code. The metrics included in the figure are, in order: the instruction schedule cost, the potential for improvement from additional machine resources, the potential for improvement from additional instruction level parallelism, the potential for improvement from vectorization, the number of retired MIAMI micro-operations<sup>1</sup>, the number of micro-operations used for address arithmetic, and the number of scalar accesses to the program stack.

Before delving into the analysis of the results, we present a short guide for the *reading and interpretation* of MIAMI results. MIAMI can output performance results in CSV format, to enable parsing with a post-processing script, and in XML format for top-down analysis using `hpcviewer`, the viewer application distributed with HPCToolkit. The screenshots included in this paper show MIAMI results displayed in `hpcviewer`. The viewer presents data in tabular format. Each metric is shown in a separate column. The rows correspond to program scopes such as loops and routines. We can expand a particular scope to understand the contribution of its inner scopes to the various metrics. The percentages visible on the right side of each cell (see Figure 3) are added by the viewer and are computed column-wise. They represent the contribution of a scope (and its children) to the total value of that metric corresponding to the entire experiment. The percentages are useful to quickly assess that a particular routine is responsible for 25% of the running time of an application, or that it accounts for 40% of the potential for improvement. However, many times we are interested in the relationship between different metrics for a particular scope. In such cases, we have to ignore the pre-computed percentages and compare the absolute values displayed inside the cells on the same row.

The results in Figure 3 are sorted by instruction schedule cost. The execution times shown in the figure do not include a memory penalty component. However, the list of top time consuming routines predicted by MIAMI matches the list and the order of routines ranked by the actual execution time measured with HPCToolkit and shown in Figure 2. Out of the three *potential for improvement* metrics computed by MIAMI, the improvement potential from vectorization is by far the largest at  $8.00e10$  out of  $1.13e11$  total instruction schedule cost. These metrics indicate that most of the code could not be auto vectorized by the compiler. While there is a large theoretical potential for improvement from vectorization, transforming the code to provide enough data parallelism in inner loops and adjusting the data layout so that the compiler can automatically vectorize a large fraction of the computation, is not a small task by any measure.

However, one metric that stands out in Figure 3 is the count of micro-operations that perform address arithmetic (metric `AddrGen`). At the entire program level,  $9.32e10$  out of  $2.75e11$  total retired micro-operations, or roughly one out of three micro-operations performs address arithmetic. However, when we look at the most time consuming routine, `ocn.time.integrator.split`, we see that this routine is responsible for 36.6% of all retired micro-operations, but it accounts for 72% of all address arithmetic instructions. In fact, just about two out of every three micro-operations executed by this routine perform address arithmetic.

Figure 4 provides a more detailed accounting of executed micro-operations broken down by functionality. We use static opcode decoding and data flow analysis to classify micro-operations. The figure includes the following metrics: the estimated instruction schedule time, the total number of retired micro-operations, the number of micro-operations performing address arithmetic, the number of micro-operations evaluating loop branch conditions, the number of scalar stack accesses, the number of micro-operations performing floating point work, and the number of micro-operations performing memory work. A few other categories are omitted for brevity.

---

<sup>1</sup>MIAMI micro-operations are very similar to x86 micro-operations, but the mapping is not necessarily 1-to-1

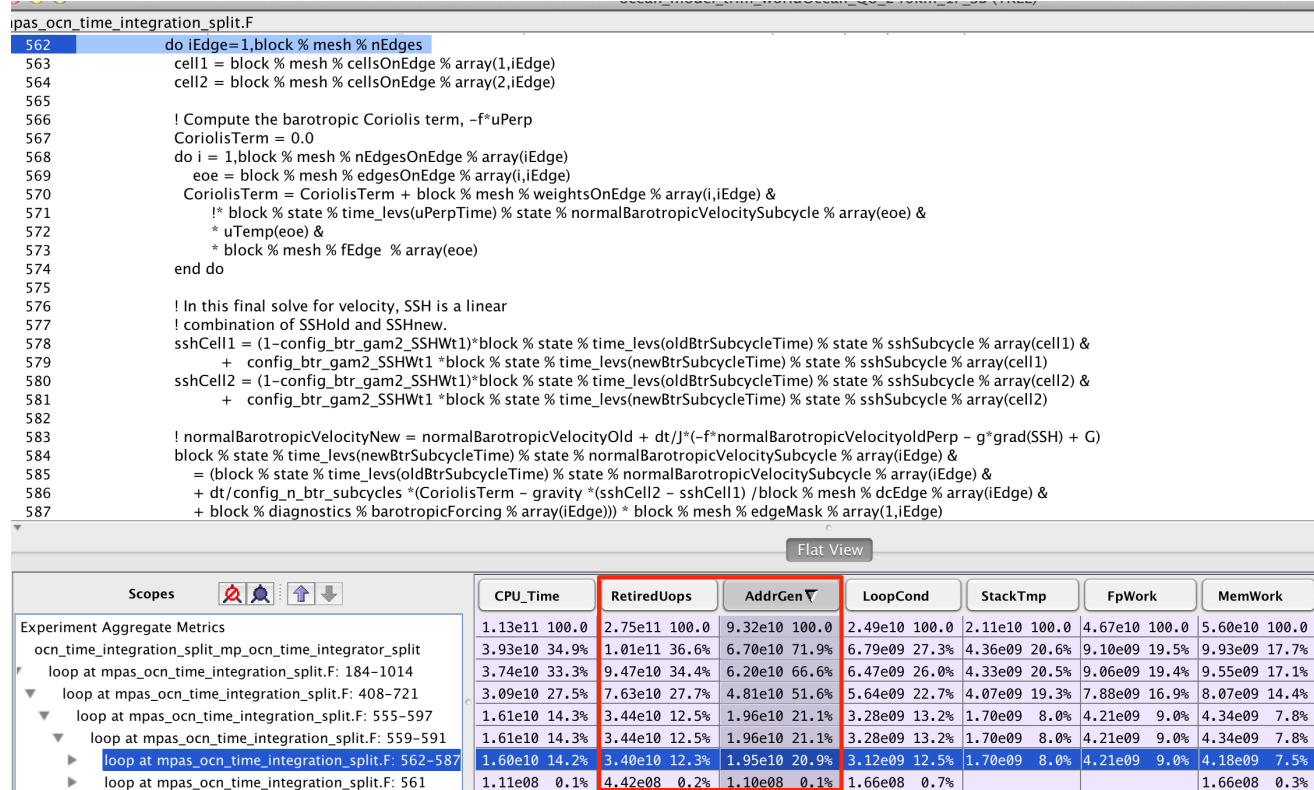


Figure 4: MIAMI performance results: breakdown of executed micro-operations by functionality.

We sorted the data by the number of address arithmetic micro-operations, and we expanded the scopes to identify the loops that contribute the highest number of address arithmetic operations. Figure 4 shows the source code of one such loop located in file `mpas_ocn_time_integration_split.F` at lines 562 - 587. There does not seem to be anything out of ordinary with this code, except the extensive use of objects of nested derived types and references to the inner components of such objects. While the compiler should be able to hoist outside the loop most of the loop-invariant address arithmetic code needed to access the inner components of these objects, our analysis shows that 57% of the micro-operations executed by this loop perform address arithmetic. We manually transformed the code to store pointers to the objects' inner components outside the loop. Inside the loop, we are using the pointers to reference the loop-dependent elements of each inner component. We verified that this approach significantly reduced execution time for the transformed loop. Next, we applied the same transformations to all the loops inside routine `ocn_time_integrator_split` that were highlighted by MIAMI as performing a disproportionately large number of address arithmetic operations.

Figure 5 shows side by side execution time profiles of the MPAS-Ocean code before and after our code transformations. The times were collected using the integrated MPAS software timers on our testbed machine. We used the 240km problem size, one MPI process, and 10 days of simulation. Figure 5a shows the costs of the main simulation phases. The time integration step clearly dominates the execution time. The transformed code achieves a 23% speedup in serial mode. Figure 5b shows the contribution of the various time integration substeps. The data preparatory substep (timer `se prep`) and the barotropic velocity prediction substep (timer `se btr vel`) benefit the most from our transformations.

Figure 6 presents a summary of the number and type of micro-operations executed by the MPAS code and by the `ocn_time_integrator_split` routine, before and after our transformations. The number

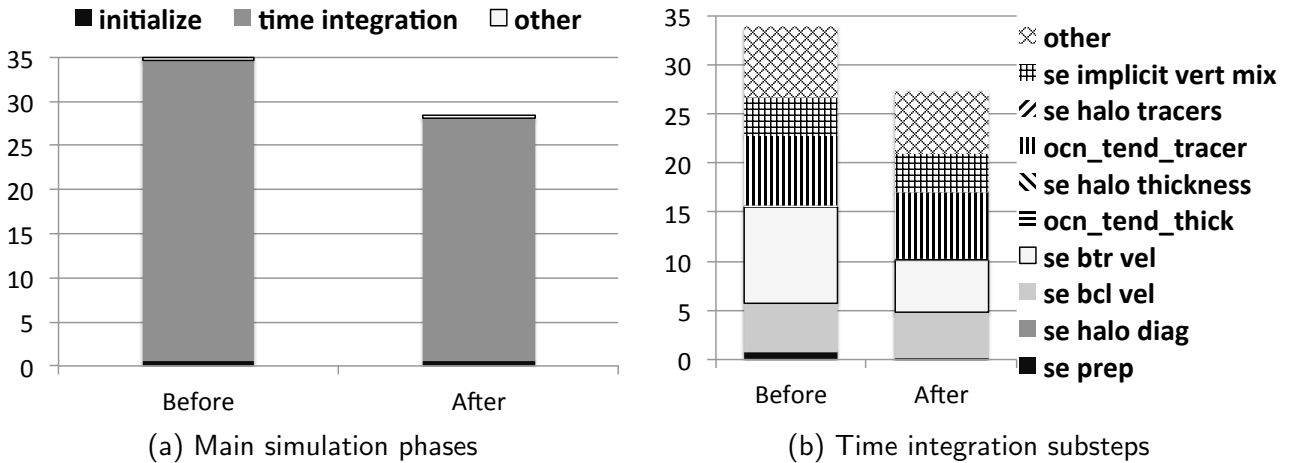


Figure 5: Execution time breakdown before and after optimizations.

	Scope	RetiredUOps	AddrGen	LoopCond	StackTemp	RegMove	FpWork	IntWork	MemWork
Before	MPAS code	2.75E+11	9.32E+10	2.49E+10	2.11E+10	8.49E+09	4.67E+10	1.67E+10	5.60E+10
	ocn_time_integrator_split	1.01E+11	6.70E+10	6.79E+09	4.36E+09	1.12E+09	9.10E+09	9.03E+08	9.93E+09
After	MPAS code	2.27E+11	4.03E+10	2.48E+10	2.59E+10	8.50E+09	4.66E+10	1.77E+10	5.58E+10
	ocn_time_integrator_split	5.28E+10	1.40E+10	6.68E+09	9.16E+09	1.13E+09	9.08E+09	2.03E+09	9.83E+09

Figure 6: Number of executed micro-operations before and after our transformations.

of address arithmetic operations dropped by almost 80% for the time integrator routine, and by 57% for the entire application. At the same time, we observe a more than doubling in the number of scalar stack references for the transformed routine. Stack references correspond to scalar temporary variables and spill / unspill code inserted by the compiler. By hoisting the bulk of the address arithmetic code outside of loops, we increased the number of concurrent live ranges in the code, which makes the register allocator’s job more difficult. However, overall, our transformations reduced the number of micro-operations executed by the `ocn_time_integrator_split` routine by 48%, which yields a 23% speedup as mentioned earlier.

There are other likely opportunities for optimization of the MPAS code’s per-core performance. However, clear understanding of these opportunities requires knowledge of the application and a careful review of the performance results produced by MIAMI and other tools.

## References

- [1] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. R. Tallent. Hpctoolkit: tools for performance analysis of optimized parallel programs <http://hpctoolkit.org>. *Concurr. Comput. : Pract. Exper.*, 22(6):685–701, April 2010.
- [2] Gabriel Marin, Jack Dongarra, and Dan Terpstra. MIAMI: A framework for application performance diagnosis. In *IEEE International Symposium on Performance Analysis of Systems and Software, 2014*. (to appear), March 2014.
- [3] MPAS: Model for Prediction Across Scales. <http://mpas-dev.github.io>.
- [4] Todd Ringler, Mark Petersen, Robert L. Higdon, Doug Jacobsen, Philip W. Jones, and Mathew Maltrud. A multi-resolution approach to global ocean modeling. *Ocean Modelling*, 69(0):211 – 232, 2013.