

NEMO: improving computational performance

ISENES3 – General Assembly – Oct, 4-6 2021



NEMO: computational performance community



Science & Technology
Facilities Council



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



Institut
**Pierre
Simon
Laplace**



Met Office



cmcc
Centro Euro-Mediterraneo
sui Cambiamenti Climatici



CERFACS

CENTRE EUROPÉEN DE RECHERCHE ET DE FORMATION AVANCÉE EN CALCUL SCIENTIFIQUE



esiwace2
CENTRE OF EXCELLENCE IN SIMULATION OF WEATHER
AND CLIMATE IN EUROPE



The IS-ENES3 project has received funding from the European Union's Horizon 2020
research and innovation programme under grant agreement No 824084

NEMO improvements

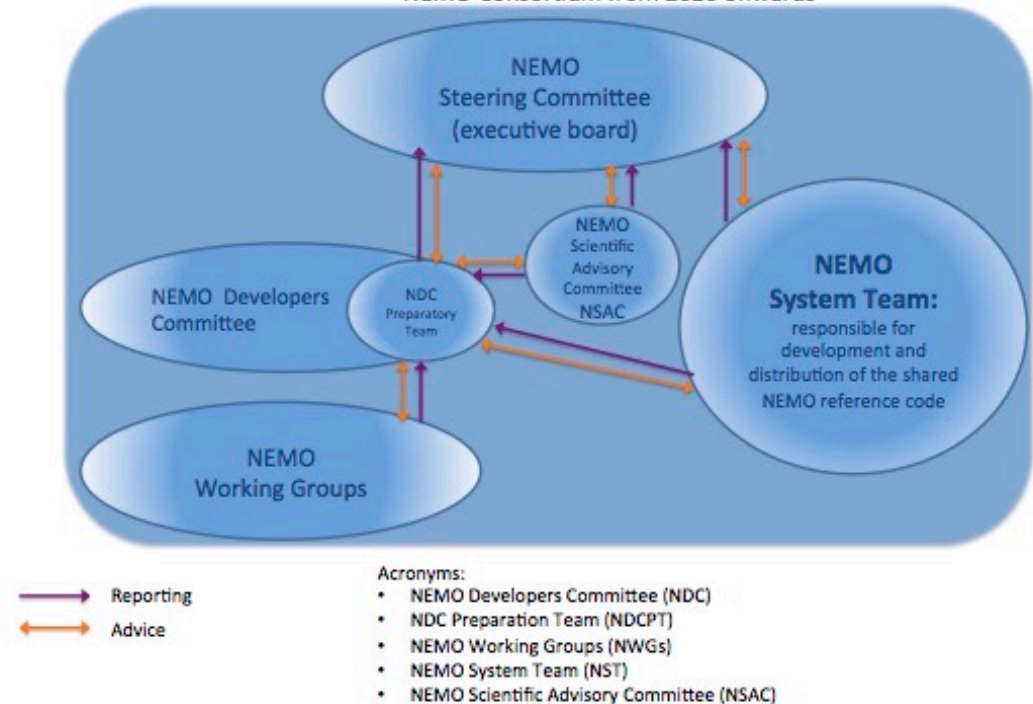
- Single core performance
 - Tiling
 - Loop fusion
 - Mixed precision
- Communication
 - Neighborhood collective communications
- Macro task parallelization
- Multigrid refinement optimization
- I/O
 - Improving read/write with XIOS
 - Online diagnostics
- Support for different architectures
 - GPU
 - DSL



NEMO Consortium organization

- NEMO System Team (NST) is responsible for development and distribution of the NEMO reference code
 - New actions are defined in the annual WorkPlan
- NEMO Working Groups articulate and coordinate the exploration of options for development of the NEMO reference code
 - NEMO HPC-WG aims at evaluating solutions to improve the computational performance of the NEMO code.

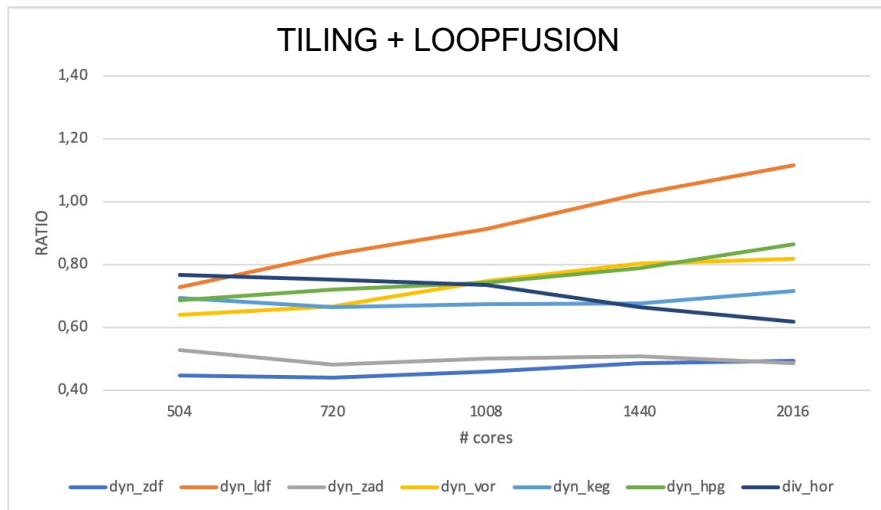
NEMO Consortium from 2020 onwards



Loop fusion and Tiling

- Efficient exploitation of memory hierarchies and hardware peak performance
- **Loop fusion technique** aims at better exploiting the cache memory by fusing DO loops together
- **Tiling** the calculation is divided into chunks of work that can remain cache-resident for as long as possible.
 - tile size and shape can be tuned appropriately for cache sizes on any platform

Loop fusion and Tiling



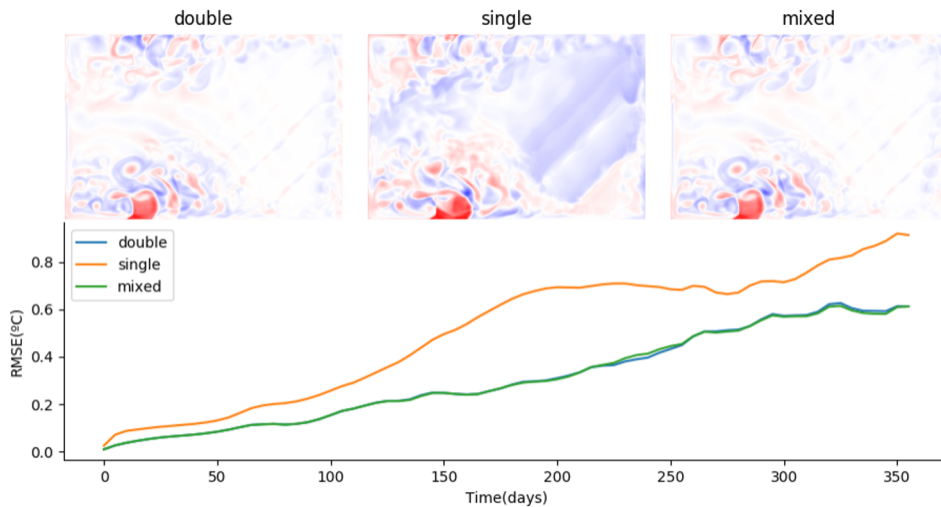
The ratio of the optimized code w.r.t. the baseline is reported changing the number of cores for the key routines of ocean dynamics. Ratio < 1 is good

- LoopFusion and Tiling applied only to the Ocean Dynamics and Ocean Tracer
- On average a speedup of 1.4x can be achieved
- The impacts of this optimization strongly depends by the platform and by the configuration



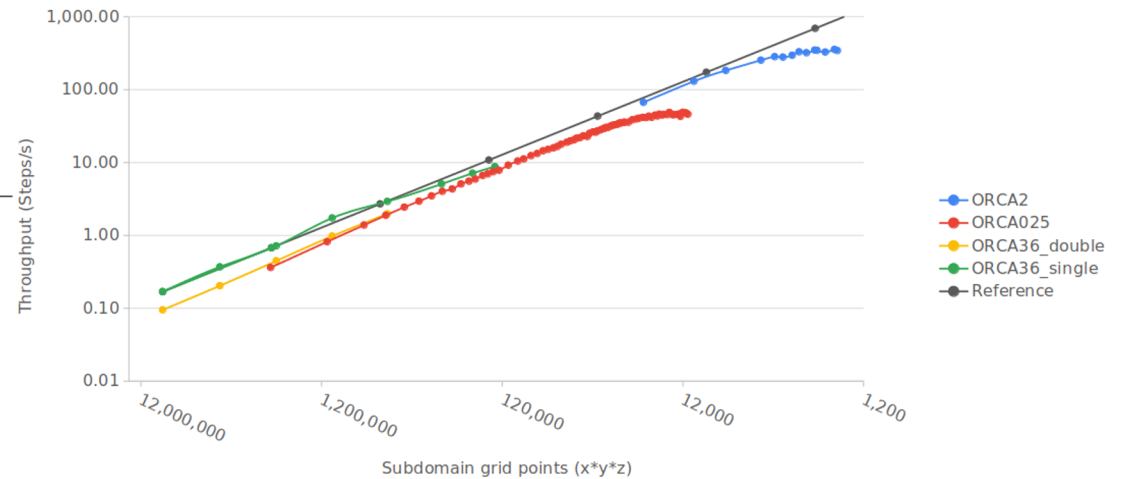
Mixed Precision

Impact of precision on sea-surface temperature in NEMO4:
comparison of GYRE1/90 simulations using different precisions

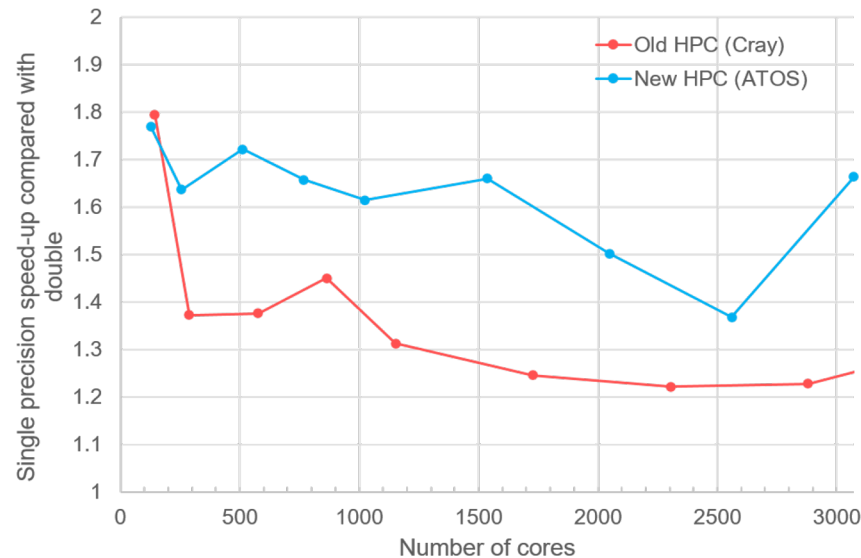


Mixed-precision approaches can provide performance benefits while keeping the accuracy of the results.

- With an appropriate tuning of the variables in SP vs those in DP, the results accuracy of the mixed precision version is preserved
- The mixed precision approach considerably improve the parallel scalability
- Mixed precision support is under evaluation to be included in the official NEMO release



Single Precision at ECMWF

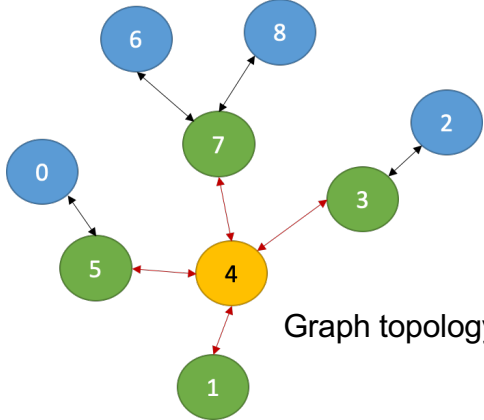
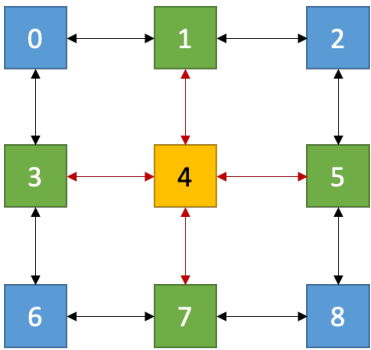


- Fully single-precision coupled atmosphere-wave-ocean forecasts now possible, **including NEMO**
- Tested with eORCA1 ocean and compared with operational reference (DP NEMO) in extended range forecasts
- Mostly skill neutral change
- Speed-up from using single precision in NEMO measured on old (Cray) and new (ATOS) HPC at ECMWF
- Final speed-up depends on I/O server → integration of NEMO with ECMWF I/O server MultIO underway



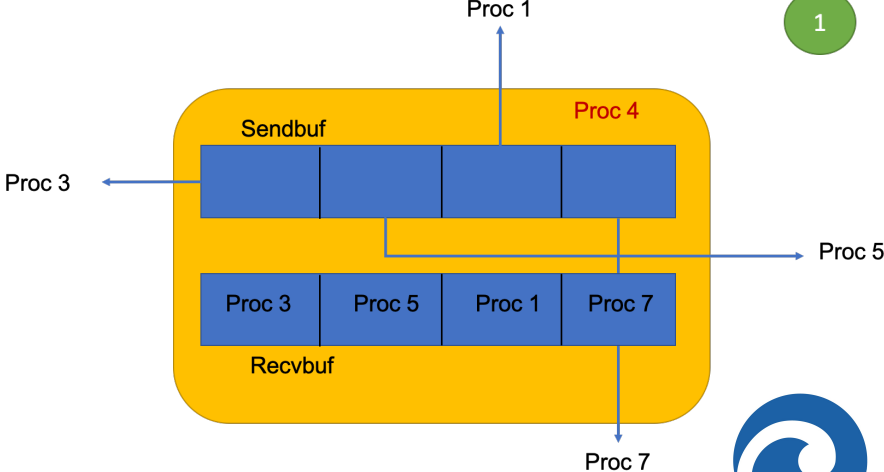
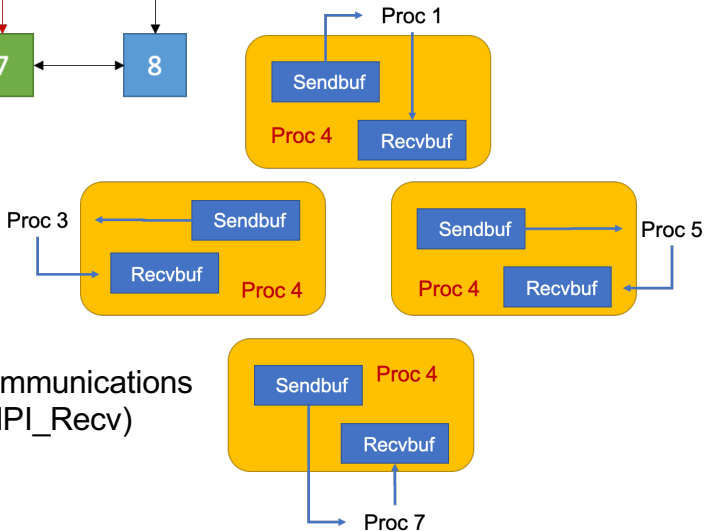
MPI Communication Neighborhood collectives

Cartesian topology



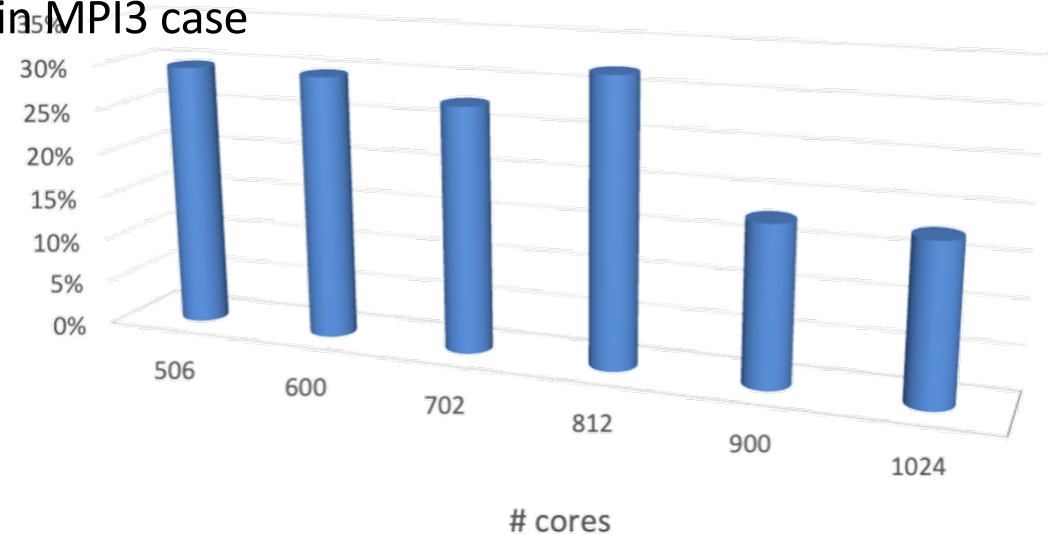
1 collective communication
(MPI_Neighbor_alltoall)

4 Point-2-Point communications
(MPI_Send/MPI_Recv)



MPI Communication Neighborhood collectives

- Extension of LBC (Lateral Boundaries Condition) module to support MPI3 Neighborhood Collectives:
 - New Cartesian communicator
 - Ranks reordering to match NEMO processes order
 - Data buffer handling
 - Implementation of multi field exchange in MPI3 case
- Test on the advection scheme
 - GYRE_PISCES configuration (nn_GYRE=200 → ~6000x4000x31 grid resolution)
 - Communication time improved within a range of 18%-32%



Macro Task Parallelism

- Parallelize OPA (ocean module) and TOP-PISCES (tracer advection biogeochemistry -BGC- module) into two executables and ensure 3D coupled fields exchange via the community coupler OASIS.
 - The ocean-BGC coupled model exhibits an improvement of computing performance when the subdomain decomposition leads to computations/communications ratio that put the performance just below the scalability limit
 - The coupling cost, caused by OASIS coupling extra cost and load imbalance between components is non negligible (around 20% in our case) but can be reduced
 - This contrasted result suggests that the only clear performance gain can only be ensured with the radical cost lowering of the most time consuming component, the BGC model (coarsening)



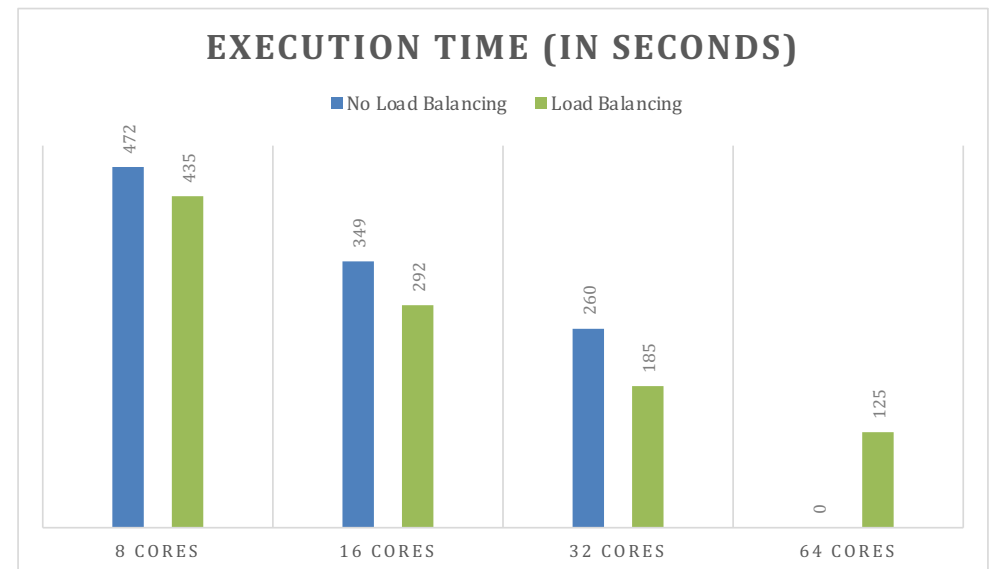
I/O optimization through XIOS

- Improvement on I/O reading initial conditions and reading regridding weights using XIOS
- the XIOS support has also been adopted for reading and writing of the restart files in the SI3 (sea ice model).



Multigrid capability

- The support for nested multigrid in NEMO is implemented in the AGRIF component
- NEMO model has been updated to provide an estimation of the computational cost of each cell grid; a new load balancing policy has been implemented in AGRIF
- Achieved 1.4x speedup on average



Inria

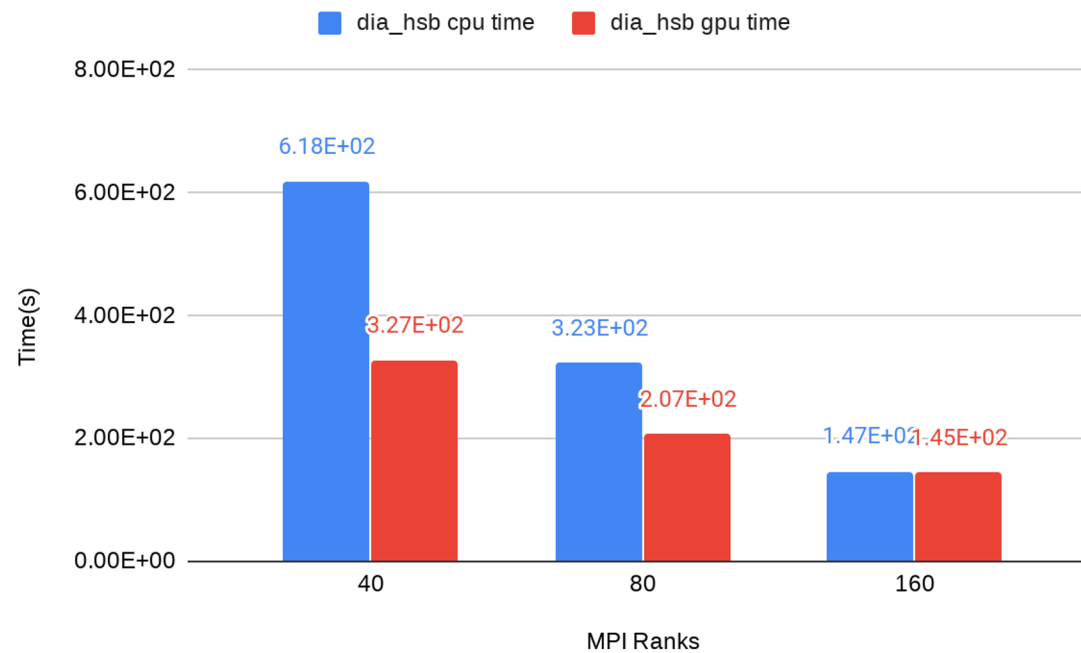


Online diagnostics – GPU based

- The rationale of this activity is to improve the NEMO computational performance by offloading the computations for diagnostics on GPU.
- The ocean global heat content, salt content and volume conservation diagnostics (`dia_hsb`) has been chosen as starting point because it is the most expensive.
- The code itself is executed 50x faster than in a single CPU but the data transfer to and from GPU is the main bottleneck.
- Pinned Memory and GPU Directly Attached to the host can be used to mitigate the data transfer penalty
- Asynchronous communications and a memory buffer approach reduce significantly the data transfer penalty

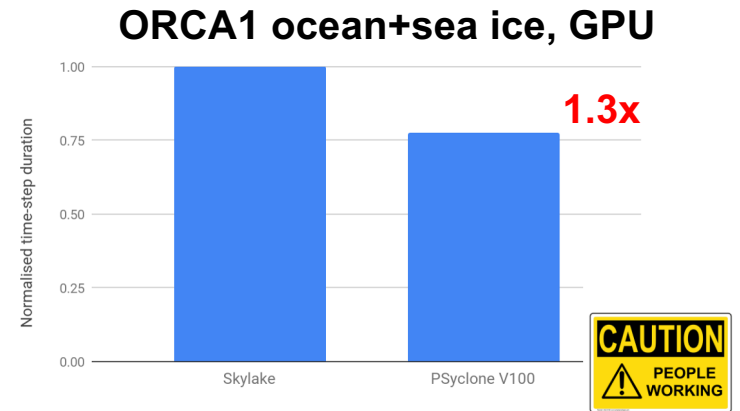
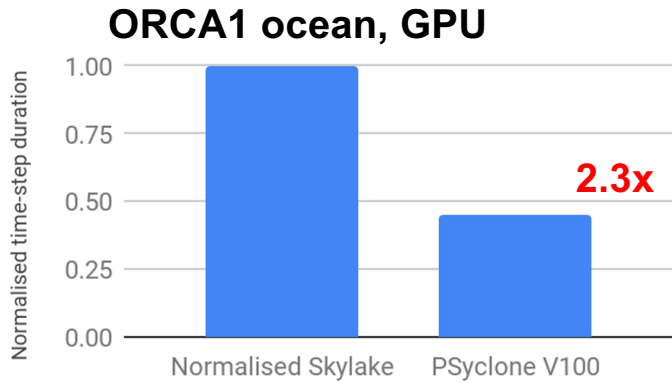
Online diagnostics – GPU based

dia_hsb scalability

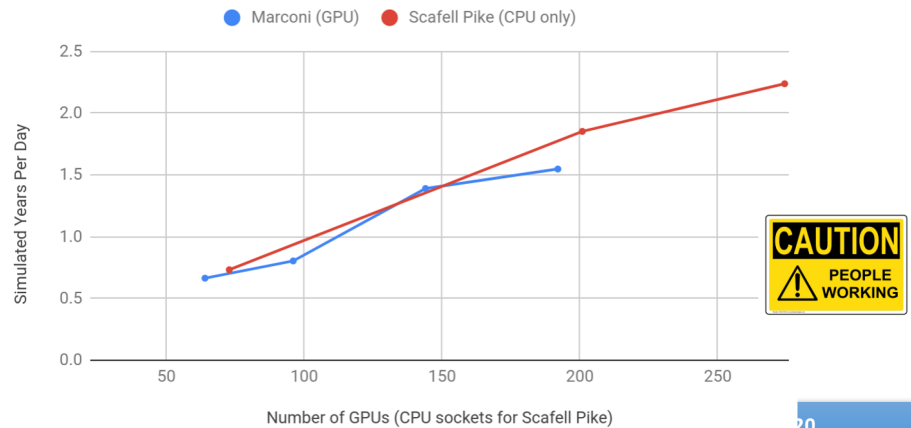


NEMO on GPU

- Use PSystem to automatically insert OpenACC directives into the code

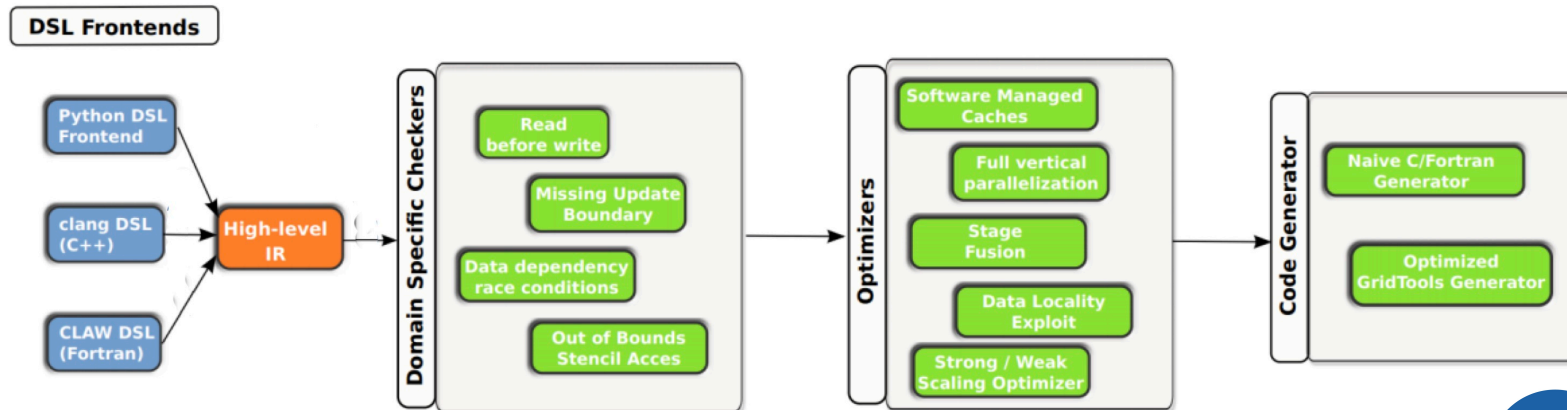


ORCA12 ocean, GPU + MPI



DSL GTClang for NEMO

- Domain Specific Language GTClang has being enhanced to support NEMO requirements (i.e. regular grid, numerical integration schema, computational kernels)
- Preliminary evaluation of GTClang through porting of specific “dwarf” which represent the advection schema (MUSCL) used in NEMO



DSL GTClang for NEMO

```

DO jk = 1, jpkm1
DO jj = 1, jpj-1
DO ji = 1, fs_jpim1
zwx(ji,jj,jk) = umask(ji,jj,jk) * ( ptb(ji+1,jj,jk,jn) - ptb(ji,jj,jk,jn) )
END DO
END DO
END DO

DO jk = 1, jpkm1
DO jj = 2, jpj-1
DO ji = 2, jpi-1
zslpx(ji,jj,jk) = zwx(ji,jj,jk) + zwx(ji-1,jj,jk)
END DO
END DO
END DO

DO jk = 1, jpkm1
DO jj = 2, jpj-2
DO ji = 2, jpi-2
zu = pun(ji,jj,jk) / ( e1u(ji,jj) * e2u(ji,jj) * fse3u(ji,jj,jk) )
zflux(ji,jj,jk) = pun(ji,jj,jk) * ( ptb(ji+1,jj,jk,jn) + zu * zslpx(ji+1,jj,jk) )
END DO
END DO
END DO

DO jk = 1, jpkm1
DO jj = 3, jpj-2
DO ji = 3, jpi-2
zu = 1. / ( e1t(ji,jj) * e2t(ji,jj) * fse3t(ji,jj,jk) )
pta(ji,jj,jk,jn) = pta(ji,jj,jk,jn) - zu * ( zflux(ji,jj,jk) - zflux(ji-1,jj,jk) )
END DO
END DO
END DO

```

```

stencil advection_MUSCL {
do {
vertical_region (k_start, k_end - 1) {
zwx = u_mask * (ptb(i+1) - ptb);
}
!-- Slopes of tracer
vertical_region (k_start, k_end - 1)
zslpx = zwx + zwx(i-1);
}
!-- Horizontal advective fluxes
vertical_region (k_start, k_end - 1) {
zu = pun / (e1u * e2u * fse3u);
zflux = pun * (ptb(i+1) + zu * zslpx(i+1));
}
!-- Tracer advective trend
vertical_region (k_start, k_end - 1) {
zu = 1.0 / (e1t * e2t * fse3t)
pta = pta - zu * (zflux - zflux(i-1));
}
}
}

```

Pros

- Easy code maintenance
- Improved code readability
- Seamless support for GPU
- Less error-prone code
- Fast and efficient technical support

Cons

- GTClang environment hard to compile and install
- No documentation
- No MPI support
- Lose of performance w.r.t. the original version

THE CONSORTIUM

Coordinated by CNRS-IPSL, the IS-ENES3 project
 gathers 22 partners in 11 countries



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement N°824084



Our website
<https://is.enes.org/>



Follow us on Twitter !
@ISENES_RI



Contact us at
is-enes@ipsl.fr



Join the community
 on ZENODO !