**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

EXCELENCIA SEVERO OCHOA

# BSC Performance tools suite: study cases on improving the efficiency of the EC-EARTH model components

Miguel Castrillo, Oriol Tintó, Kim Serradell
Francisco J. Doblas-Reyes, Jesús Labarta

# Efficiency in Earth science models

- Efficiency is especially **critical** for Earth science models

- Simulations use a **huge amount** of computational resources

- Future simulations will need **many more resources**
  - Computational time
  - Storage and postprocess
  - Software to simplify the usage of the model
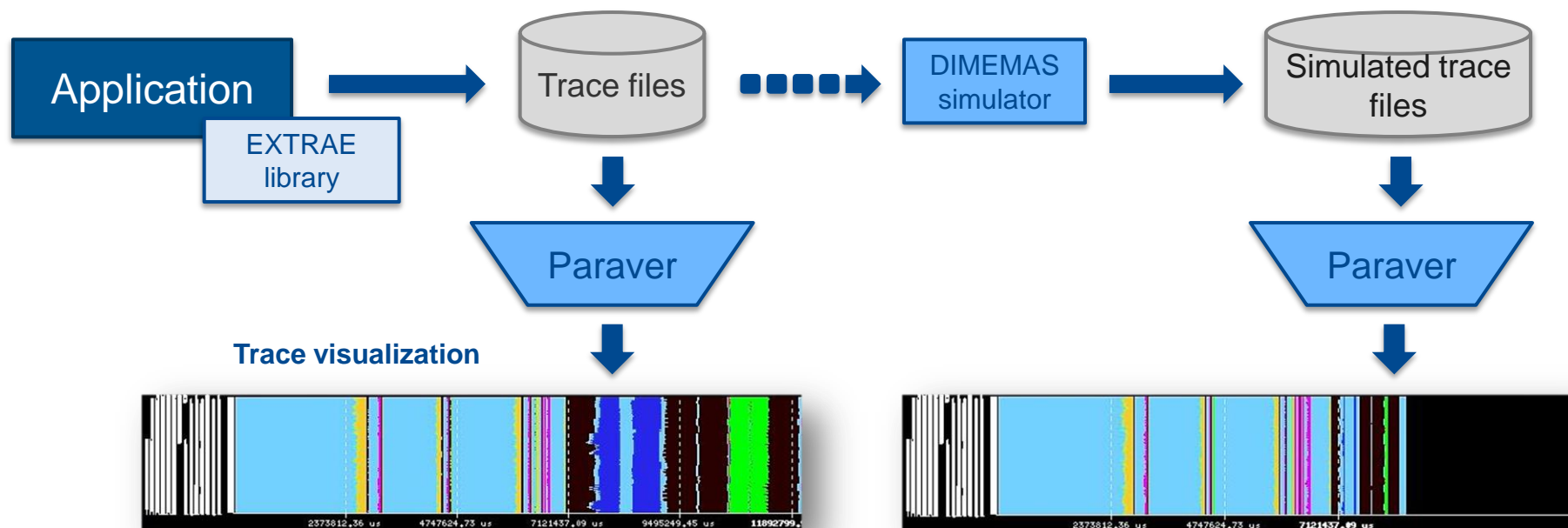
# Energy efficiency

- Energy efficiency

**Increase performance**

$$EE = \frac{Performance}{Power\ consumed}$$

Reduce power

- Performance loss caused by:
  - Bad programming
  - Load imbalance
  - Synchronization
  - Resource contention
  - …

# BSC Computer Sciences Performance Tools

- Since 1991
- Based on traces
- Open Source: http://www.bsc.es/paraver
- **Extrae**: Package that generates Paraver trace-files for a post-mortem analysis
- **Paraver**: Trace visualization and analysis browser
  - Includes trace manipulation: Filter, cut traces
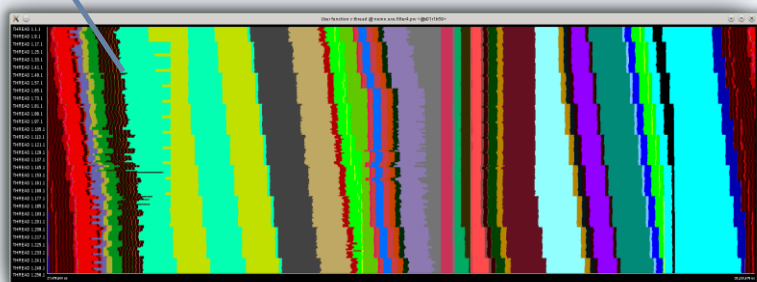- **Dimemas**: Message passing simulator



**Trace visualization**

DIMEMAS generated trace. Target = ideal machine

# PARAVER trace analysis

## Serial efficiency

**Functions instrumentation**

Each color represents a function



Useful duration



Darker color = Longer duration

Useful IPC - Instructions per cycle



Lighter color = Less IPC

Horizontal axis -> Time component

## Correlation between two functions



Darker color = Higher values

## Parallel efficiency



MPI stats reflect the percentage of time invested in computation for each thread.

Total stats give the communication efficiency and the load balance

\* Examples from NEMO 3.4–LIM2 / ORCA025

4

# EC-Earth: A coupled climate model

- Earth System Model

- Reliable in-house predictions of global climate change

- Part of a Europe-wide consortium

- Being used in **large** European **projects**
  - EMBRACE
  - EUPORIAS
  - IS-ENES
  - SPECS
  - PRIMAVERA

- 3.1 version → IFS + NEMO-LIM + OASIS
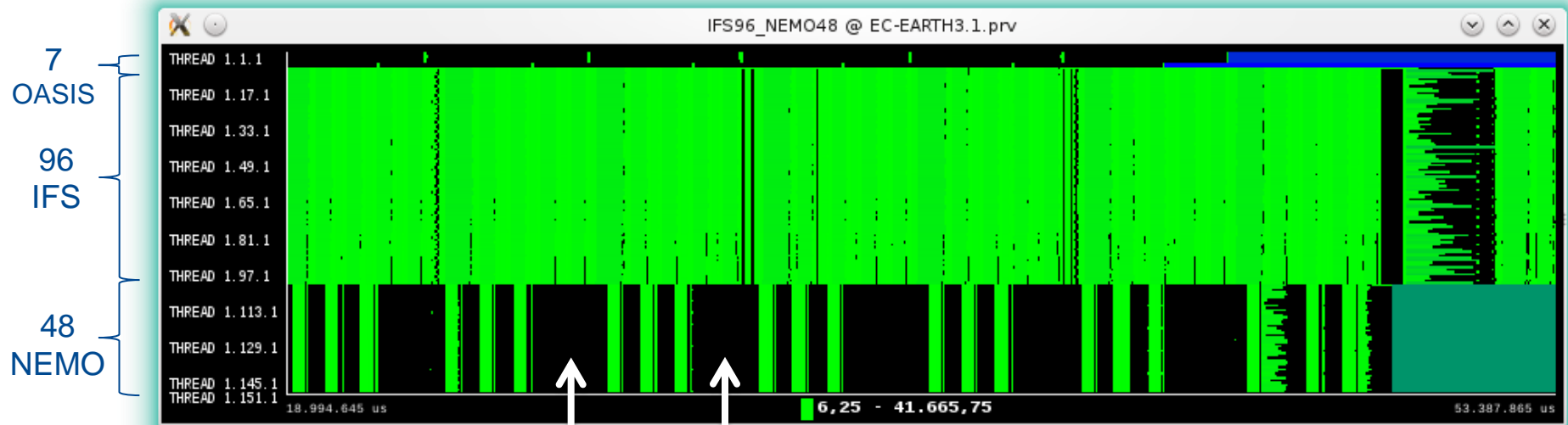
# PRIMAVERA Project

- Preparing the next CMIP6 high-resolution simulations called HiResMIP
  - Target resolutions are T511-ORCA025 and T1279-ORCA12

- 19 European groups involved

- No experience in analyzing efficiency on these resolutions

| Resolution | Single climate experiment (10 members, 60 start dates, 10 years simulated) | | |
| --- | --- | --- | --- |
| | Grid Size | Output Size | Computation Time |
| T511-ORCA025 | Atmos 40km - Ocean 25km | 720 Tb | $132.0 \times 10^6$ |
| T1279-ORCA12 | Atmos 16km - Ocean 9km | 1,4 Pb | NA |

# An EC-Earth Paraver trace

- Motivation: Finding a good configuration to **optimize** the resources usage
- IFS T255L91-ORCA1L46
- Configuration widely used in **production**
  - Using 7 cores for OASIS, 96 for IFS and 48 for NEMO
- 1 day simulation traces
- Traces generated in burst mode (only computational regions > 100us)
- Paraver view → Useful duration (displays duration of computational bursts)



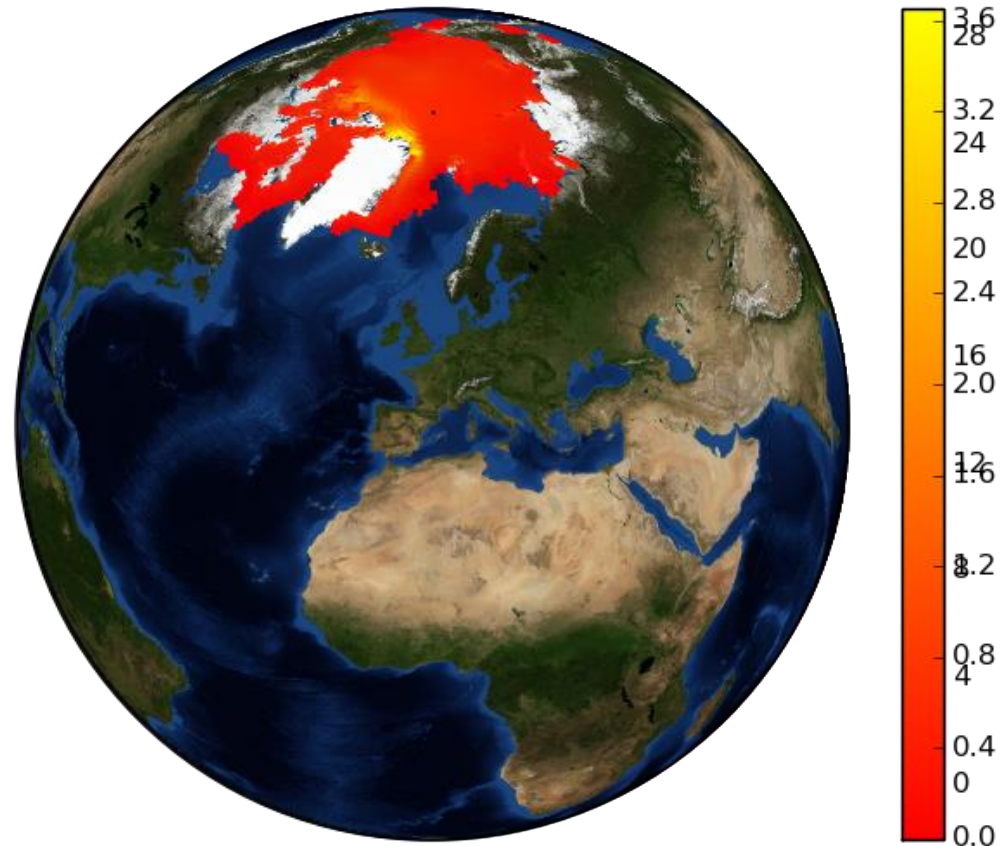Black regions → Not computation (MPI, I/O…) → NEMO waiting

*Time axis*

# NEMO: An ocean model

**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

EXCELENCIA
SEVERO
OCHOA



**Nucleus for European Modeling of the Ocean (NEMO)** is a **state-of-the-art** global **ocean model**

It is used in oceanographic research, operational oceanography, seasonal forecast and climate studies
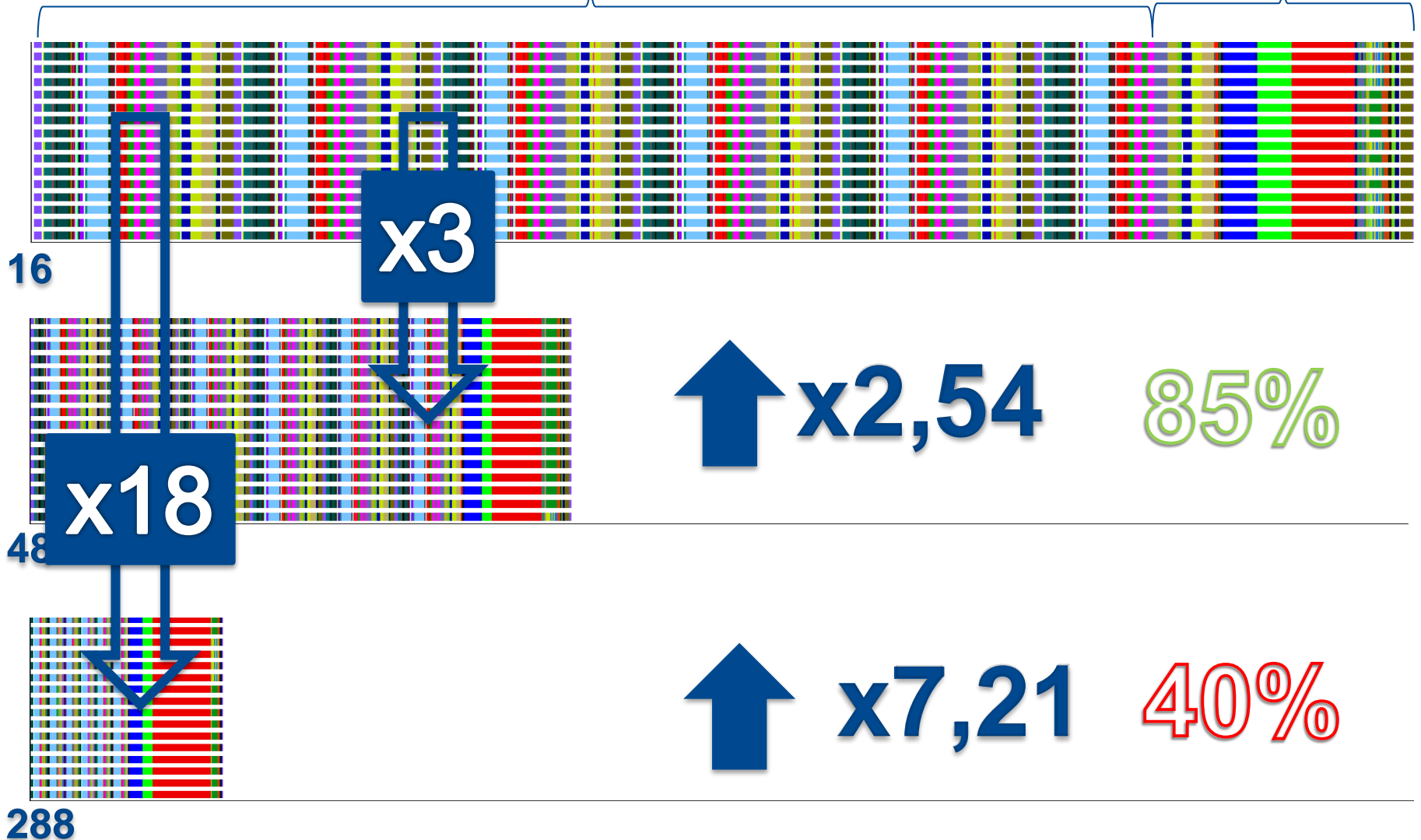
Includes several **sub-models**. Many of them can work in standalone version , many others need to be coupled
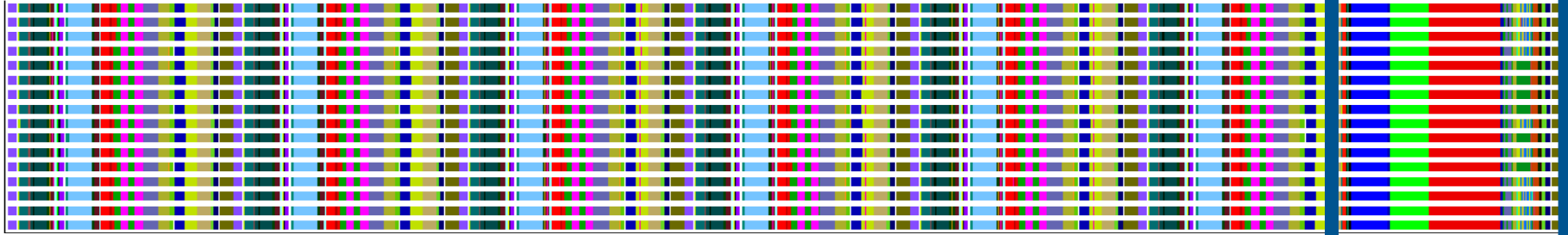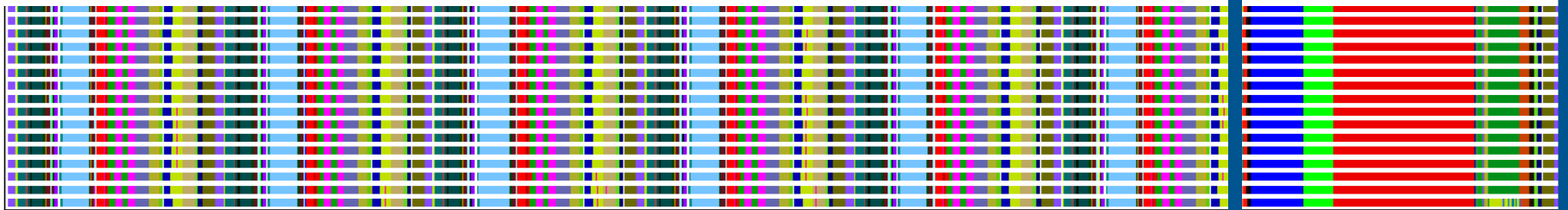


Sea Ice Thickness
Sea Surface Temperature

**OPA**   **LIM**

x3

x18

16

48

288

x2,54   85%

x7,21   40%

*Timelines have the same duration

**16**

**48**

**288**

## LIM HDF

**288**



Outside MPI
MPI Isend
MPI Recv
MPI Wait
MPI All_gather
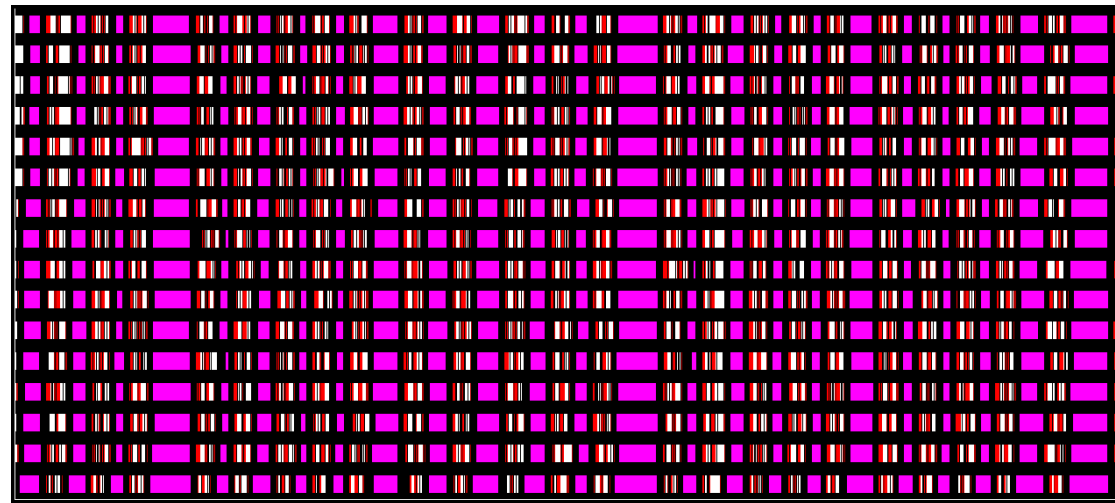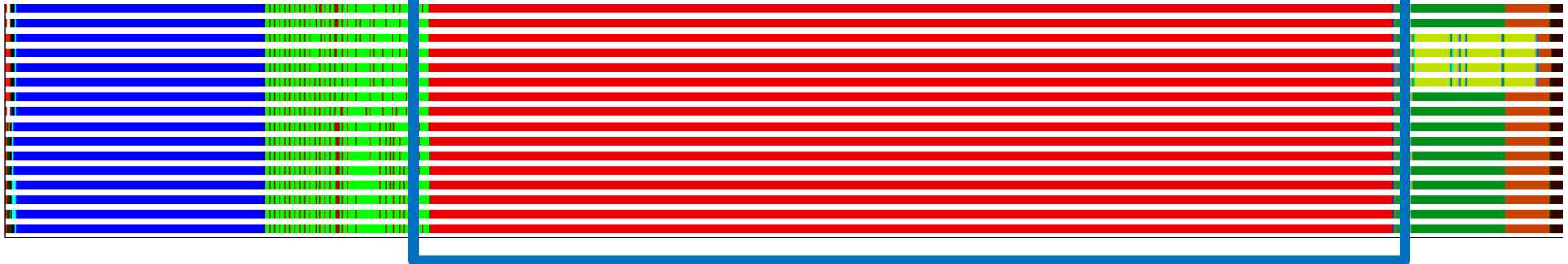
Only **20%** of the time invested on **computation**

Global Communication at **every** loop **iteration → 60%** of the time

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

EXCELENCIA
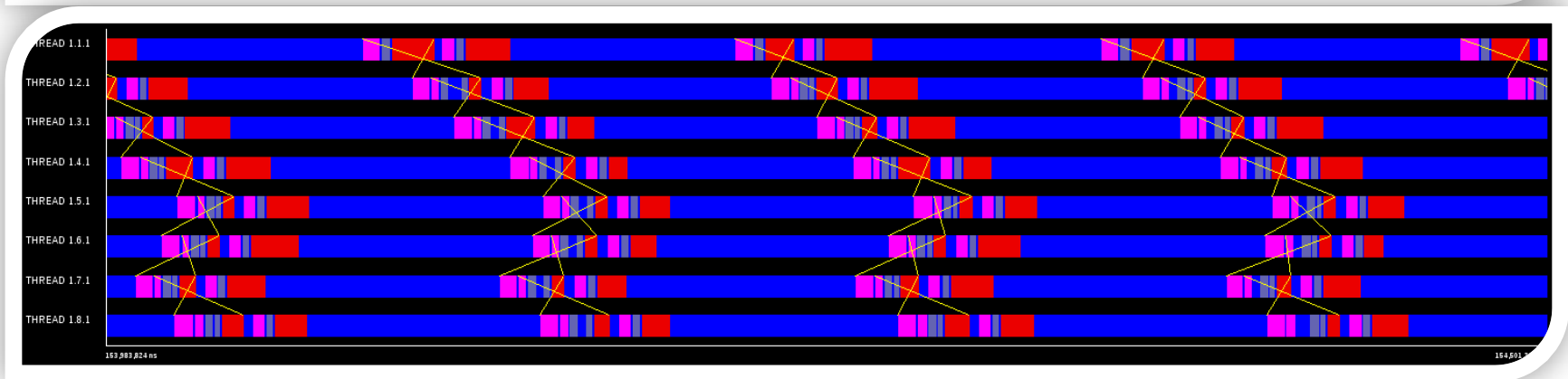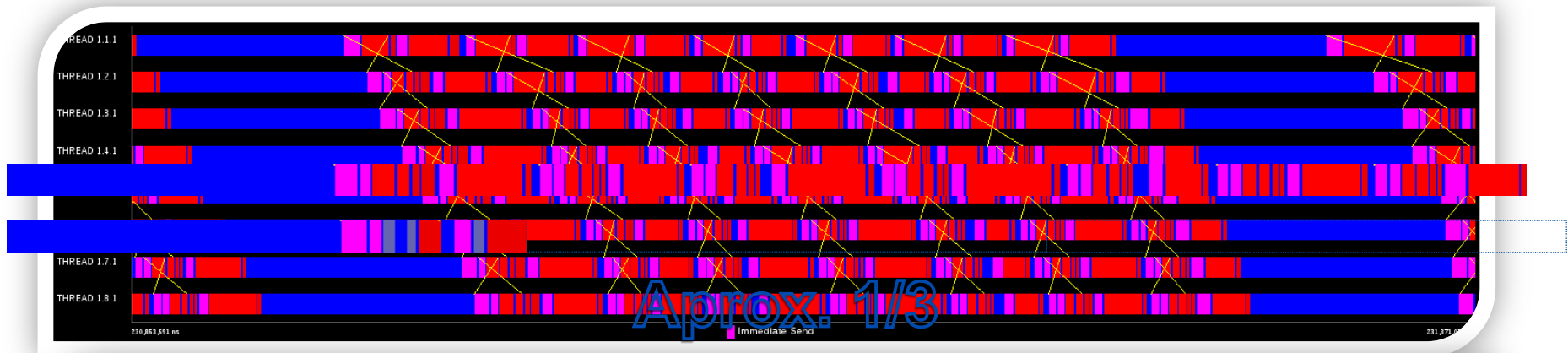SEVERO
OCHOA

# NEMO model optimizations

- Convergence check that is executed in every loop iteration
- Control structure put to reduce the frequency by a factor of N

```fortran
do while( control > threshold)    ! Sub-time step loop
    ...
       some computation with x
    ...
    !
    call interchange ( x )  !  lateral boundary condition
    !



    control = max( x )  ! Find local max
    !
    call global_max( control )   ! Find global max
    !



end do                          ! end of sub-time step loop
```
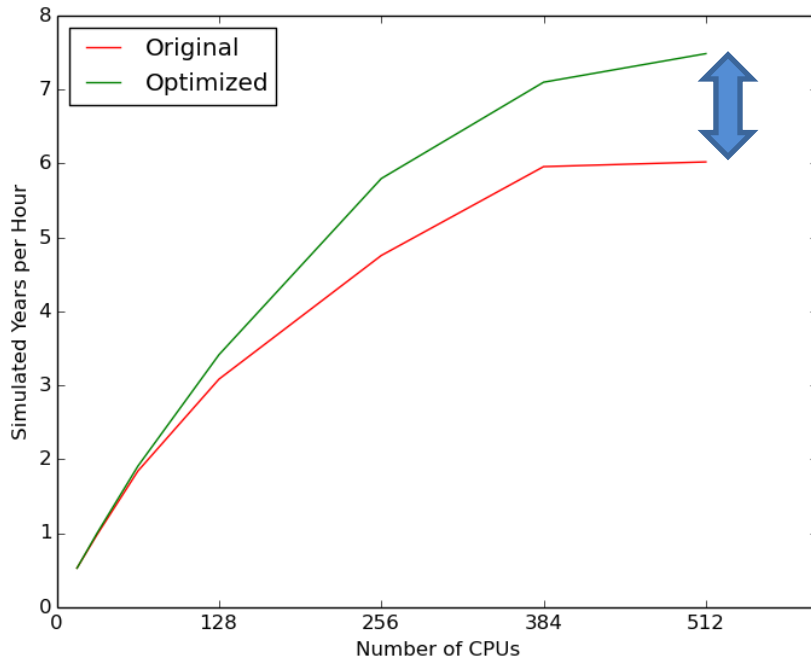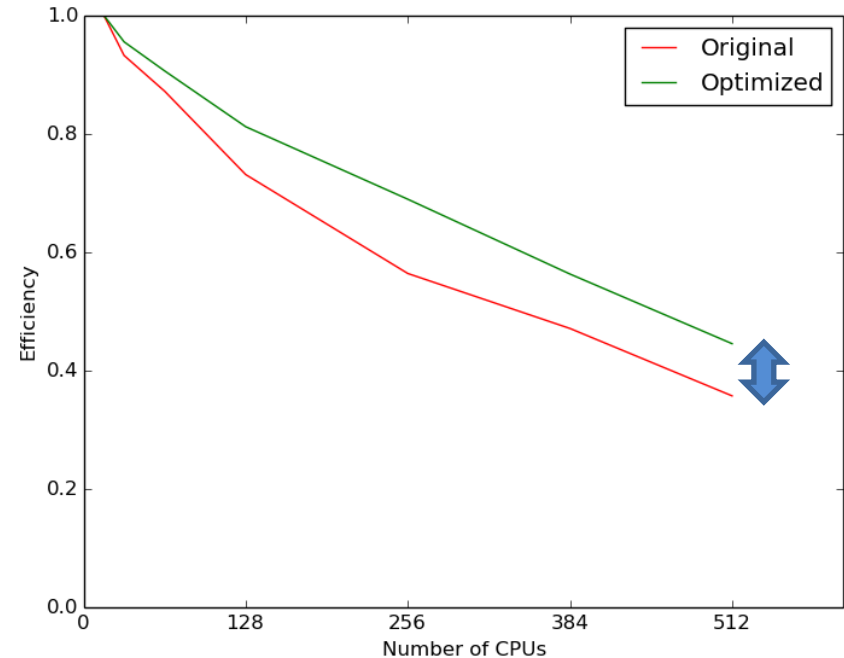
# Message packing

# Performance improvement

- Message Packing + Reduction of global communications
- Increase in the model scalability and efficiency



**Scalability Improved!**     **Efficiency Improved!**

# But it can be even better…

- Even though, in each LIM iteration we have:
  - 41 **lim_hdf** calls
  - More than 1400 collectives and border interchanges
- **lim_hdf** calls are *(almost)* independent → Reorder it to achieve coarser granularity and reduce collectives number by using the message packing

**Original**



9390 Isend/Recv – 1163 All_Reduce

**Message packing + Convergence check**



9659 Isend/Recv – 393 All_Reduce

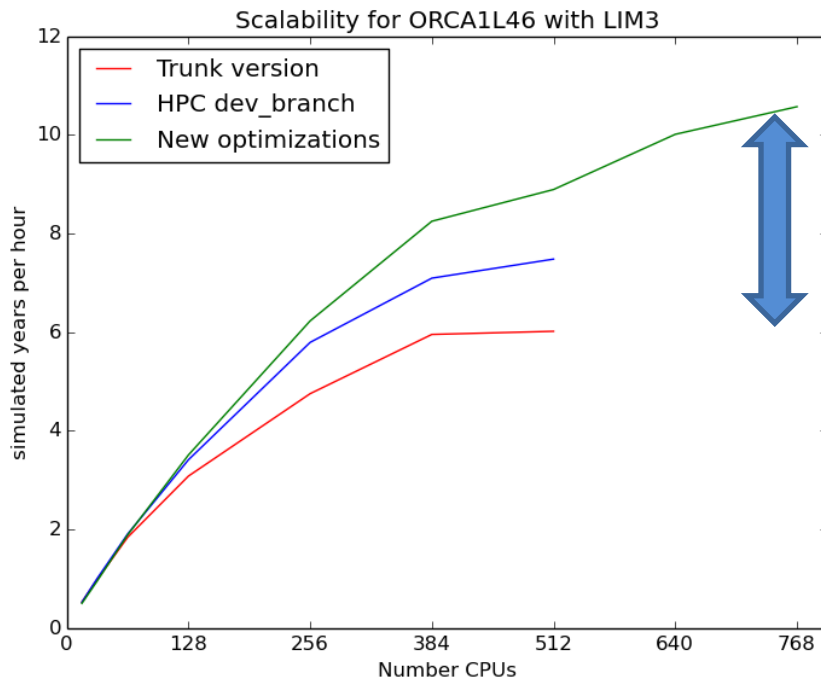**Message packing + Convergence check + lim_hdf multiple**



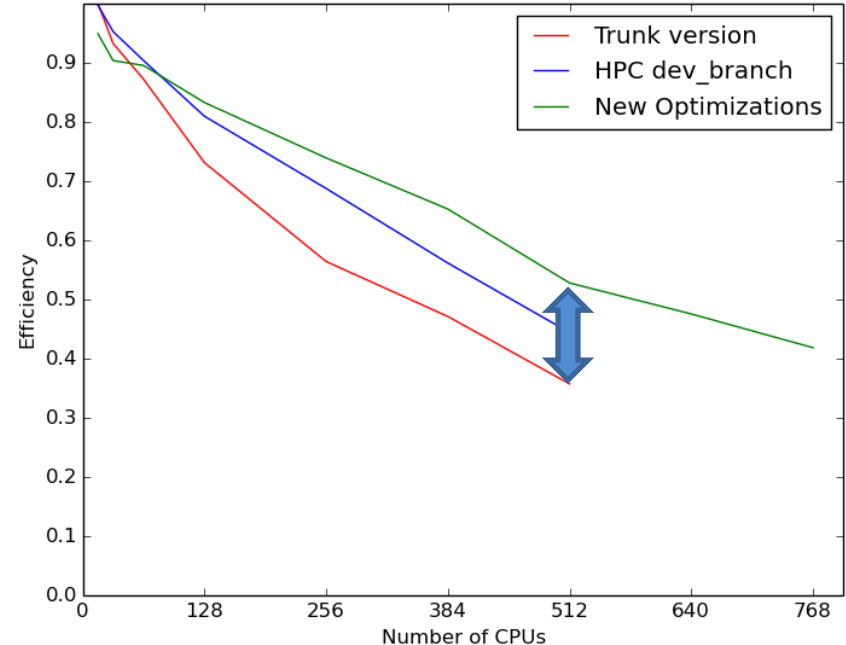**6432** Isend/Recv – **163** All_Reduce

-32%!          -86%!          **Just reordering one routine!**

# More performance improvement

- Previous optimizations already included in a new model branch and now are merged into the NEMO 3.6 trunk
- New optimizations increase further scalability



Scalability for ORCA1L46 with LIM3

**Scalability Improved!**    **Efficiency Improved!**

Conclusions

# Conclusions

- **Little changes** in the configuration can significantly improve the performance

- **Trace analysis** can **guide** the **users** in understanding the behavior of the code

- A precise analysis and prediction can generate ideas that **direct** the **restructuring** of the application in the most productive way

- Performance analysis is critical for a **rational usage** of the **resources**

# Thank you!

For further information please contact
miguel.castrillo@bsc.es