

Barcelona Supercomputing Center Centro Nacional de Supercomputación





Computational Profiling Analysis for Climate and Weather

Mario C. Acosta and Miguel Castrillo

26/08/2021

Summer School on Effective HPC

# Computational Profiling Analysis for Climate and Weather

### Objectives

- Define performance analysis fundamentals (objectives, methods, metrics, hardware counters, etc.)
- Define a methodology to study HPC performance for numerical models, know your enemy.
- Describe the BSC performance analysis tools suite (Extrae, Paraver, Dimemas)
- Interpret uses cases from Earth System Models (HARMONIE, IFS, NEMO, etc.) that illustrate how to identify and solve performance issues
- Apply profiling techniques to identify performance bottlenecks in your code
- Summarise typical performance problems
- Discuss specific knowledge about performance analysis applied to earth system modelling













#### Accuracy of PMSL forecast (in days) compared to baseline of 1-day forecast in 1980



**NCEP Operational Forecast Skill** 



















#### Performance Development

Lists















- To be able to use the computing power of modern supercomputers, applications must exploit parallelism.
- Parallelism produce overhead (extra computation and communications)
  - "Overhead does not look a problem in my model" → But if the needs increase (i.e. higher resolutions), a bad implementation will be a problem in some point.
  - We need a method to evaluate the parallelism efficiency of our computational models.
    - When the hardware change
    - When the number of resources change
    - When the model complexity change
    - When the resolution change
    - ...















- The necessary refactoring of numerical codes is given a lot of attention and is stirring a number of discussions.
  - Computational performance analysis and new optimizations are needed for actual numerical models.
  - Study new algorithms for the new generation of high performance platforms (path to exascale).
- Several European institutions and projects working together on the same direction (ESCAPE2, ESiWACE2, IS-ENES3, ETP4HPC...)



AND CLIMATE IN EUROP

## **CES-Performance Team & ES Department**



- Knowledge about the mathematical and computational side of Earth System Applications
- Knowledge about the specific needs in HPC of the Earth System Applications
- Researching about HPC methods specifically used for Earth System Applications



**Supercomputing** Supercomputing Center Centro Nacional de Supercomputación













- Mathematical study
  - Some methods could be better than others
    - Discretization used (explicit, implicit, semi-implicit...)
    - Parallel adaptation (solvers, preconditioners...)
  - How to implement new algorithms for new architectures
- Computational study
  - Achieve load balance among components
  - Reduce overhead introduced by parallel applications
  - Assure that the computational algorithm takes advantage of the architecture



Start Computational and mathematical study Scalability Study Performance analysis (profiling and/or tracing) Optimization **Results** verification Execution time measurement ls it good No enough l, Yes Integrate optimization S overall No speedup good? , Yes

End

- 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





End

Introducing optimizations

----

- Improvement of the mathematical and/or computational algorithm
  - Apply scientific methods which are found in the literature
  - Improve the method with a new approach
  - Revolution: Create a new (and better)
     algorithm taking into account the research line
     followed



Reproducibility study

•

٠

- Evaluate if the accuracy and reproducibility of the model is similar using or not the optimizations proposed
- Take into account the nature of climate models
  - How to evaluate, in parallel executions, if the differences





- Reproducibility study
  - Evaluate if the accuracy and reproducibility of the model is similar using or not the optimizations proposed
  - Take into account the nature of climate models
    - How to evaluate, in parallel executions, if the differences between runs are significant or not.

Kolmogorov-Smirnov differences of two 5-members ensambles





## **Profiling Analysis: BSC Tools**

- BSC Tools
  - General description
  - Extrae
    - General description
    - How to use it
  - Paraver
    - General description
    - How to use it
    - Configurations available
  - Dimemas
    - General description
  - How to work with large traces
    - Filtering/Burst mode
    - Cutting

### **BSC Tools**

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



### **BSC** Tools

### CORE TOOLS



### https://tools.bsc.es/downloads

### **BSC** Tools

### Home » Documentation » Tutorial Guidelines

These six tutorials can be opened with wxParaver versions newer than 4.3.0, and you'll be able to follow the steps within the tool. To install them, download and untar the package and follow the instructions of the Help/Tutorial option on the Paraver main window. You can download them in a single package either in .tar.gz format (127 Mb) or .zip format (127 Mb).

- Paraver introduction (MPI) Start here to familiarice with Paraver basic commands and the first steps of a performance analysis.
- Dimemas introduction The basic steps to learn how to configure and run the Dimemas simulator and to start looking at the results.
- Introduction to Paraver and Dimemas methodology This tutorial presents different ways to analyze a MPI application through well-known rules, their diagnosis and how they impact on your exploration (no traces included).
- Methodology This tutorial shows some examples of the analysis that can be done using the provided configuration files.
- Tutorial on HydroC analysis (MPI, Dimemas, CUDA) One example of performance analysis of the MPI application Hydro and further simulations with Dimemas.
- Trace preparation Look at this tutorial to select a representative region for a large trace that cannot be loaded into memory.
- Trace alignment tutorial. If you identify some unexpected unnalignement or backwards communications, use this tutorial to learn how to correct shifts between processors.

### Methodology of analysis

MPI+OpenMP Performance Analysis tips

#### **Tutorial slides**

Introduction

Core tools	Advanced features
Paraver, Detailed material	Tools scalability
Dimemas	Clustering
Extrae	Sampling

### https://tools.bsc.es/tutorial\_guidelines

- Trace generation





Barcelona Supercomputing Center Centro Nacional de Supercomputación -----





- Trace Generation: Set Environment
  - Module load extrae
    - load extrae 3.X.0 (PATH, EXTRAE\_DIR, EXTRAE\_ROOT, EXTRAE\_LIB)
  - Job script  $\rightarrow$  trace-fortran.sh | trace-c.sh
  - Extrae config  $\rightarrow$  extraeMPI.xml | extraeMPI+OMP.xml
  - Files modified for model  $\rightarrow$  run\_parallel.sh







- Job script:trace-fortran.sh
  - Available loading extrae module

```
#!/bin/bash
#Workaround for tracing in MN3, make TMPDIR point to an existing dir
#if [ ! -z "${TMPDIR}" ]; then
        export TMPDIR=$TMPDIR/extrae
#
#
        mkdir -p $TMPDIR
#fi
EXTRAE_ROOT=/usr/local/apps/extrae
if [ -z "$EXTRAE_DIR" ]; then
    echo "ERROR: EXTRAE_DIR not set, maybe extrae module not loaded?"
    exit -1
                                 True if openmp is used
fi
if [ -z "$OMP_TRACE" ; then
                                                                  Extrae config by default
    export LD_PRELOAD=${EXTRAE_DIR}/lib/libmpitracef.so
    if [ -z "$EXTRAE_CONFIG_FILE" ]; then
        export EXTRAE_CONFIG_FILE=${EXTRAE_ROOT}/xml/MPL/extrae.xml
    fi
else
    export LD_PRELOAD=${EXTRAE_DIR}/lib/libompitracef.so # For Fortran apps
    if [ -z "$EXTRAE_CONFIG_FILE" ]; then
        export EXTRAE_CONFIG_FILE=${EXTRAE_ROOT}/xml/MPI+OMP/extrae.xml
    fi
fi
```

- Extrae config:extrae.xml
  - Available using nama\_CY43R1\_IFS\_traces branch



- Extrae config:extrae.xml
  - Available using nama\_CY43R1\_IFS\_traces branch

```
<storage enabled="no">
   <trace-prefix enabled="yes">TRACE</trace-prefix>
   <size enabled="no">5</size>
   <temporal-directory enabled="yes">/scratch</temporal-directory>
   <final-directory enabled="yes">/gpfs/scratch/bsc41/bsc41273</final-directory>
 </storage>
 <buffer enabled="yes">
   <size enabled="yes">500000</size>
   <circular enabled="no" />
 </buffer>
 <trace-control enabled="no">
   <file enabled="no" frequency="5M">/gpfs/scratch/bsc41/bsc41273/control</file>
   <global-ops enabled="no"></global-ops>
   <remote-control enabled="no">
     <signal enabled="no" which="USR1"/>
   </remote-control>
 </trace-control>
 <others enabled="yes">
   <minimum-time enabled="no">10M</minimum-time>
   <finalize-on-signal enabled="yes"
     SIGUSR1="no" SIGUSR2="no" SIGINT="yes"
     SIGQUIT="yes" SIGTERM="yes" SIGXCPU="yes"
                                                   Emit computation burst of a minimal duration
     SIGFPE="yes" SIGSEGV="yes" SIGABRT="yes"
   1>
   <flush-sampling-buffer-at-instrumentation-point enabled="yes" />
 </others
                                                    Plus summarized MPI events
  bursts enabled="no">
   <threshold enabled="yes">500u</threshold
   <mpi-statistics enabled="yes
   bursts>
 <sampling enabled="no" type="default" period="50m" variability="10m" />
 <dynamic-memory enabled="no">
   <alloc enabled="yes" threshold="32768" />
   <free enabled="yes" />
 </dynamic-memory>
                                       Merge individual traces automatically
 <input
               enabled="no" />
  <merge enabled="yes"
                     efault"
    synchronization=
   tree-Tan-out="16"
   max-memory="32768"
   joint-states="yes"
   keep-mpits="yes"
   sort-addresses="ves"
   overwrite="ves"
 1>
</trace>
```

- Files modified for run\_parallel
  - If BSCTRACE=1  $\rightarrow$  Extrae is used

```
if [[ %BSCTRACE:0% = 1 ]]; then
      command="aprun -cc cpu $marg -n $submit_total_tasks -N $submit_tasks_per_node $Sarg -j $submit_cpus_per_compute_unit -d $omp_num_threads $submit_force_numa_memory_afficity trace-fortran.sh $(wher
                                                                                                                                         $cm
d) $args"
     else
      command="aprun -cc cpu $marg -n $submit_total_tasks -N $submit_tasks_per_node $Sarg -j $submit_cpus_per_compute_unit -d $omp_num_threads $submit_force_numa_memory_affinity $(whence $cmd) $args"
    fi
   fi
                                                                 Parameter to activate profiling
                                        IBRARY_PATH=${PIFS_LD_LIBRARY_PATH:-$LIBS}
                     export EC
                     if [[%BSCTRACE:0% = 1]]; then
                        LOAD_MODULE extrae
                        export EXTRAE_CONFIG_FILE=/perm/rd/nama/extrae/xml/MPI/extrae.xml
                        export TRACEDIR=${WDIR}/bsctrace.${TASK}.${ECF_TRYNO}.$$
                        export EC_LD_LIBRARY_PATH=${EC_LD_LIBRARY_PATH}:$EXTRAE_LIB
                       if [[ %OMPTRACE:0% = 1 ]]; then
                         export OMP_TRACE=1
                        fi
                       mkdir $TRACEDIR
                     fi
                     (export LD_LIBRARY_PATH=$EC_LD_LIBRARY_PATH ; eval $command)
                     if [[ %BSCTRACE:0% = 1 ]]; then
                                                                       Trace files generated
                            .prv *.pcf *.row STRACEDIR
                        mv
                     fi
```

### **BSC Tools:**Paraver



Comparative analyses Multiple traces Synchronize scales

### **BSC** Tools:Paraver

Barcelona Supercomputing Center

- **Paraver traces:** made up from records (timestamp + event or activity) of three different kind:
  - **State records:** intervals of thread status, i.e., waiting in a barrier (either MPI or OpenMP), waiting for a message, computing...
  - **Event records:** punctual event occurred in a given timestamp, as entry & exit points of user functions, MPI routines, OpenMP parallel regions...
  - **Communication records:** relationship between two objects, as communication between two processes (MPI), task movement among threads (OpenMP/OmpSs) or memory transfers (CUDA/OpenCL).





## How a trace looks like: basic overview



AND CLIMATE IN FURDER

entro Nacional de Supercomputación

## **BSC Tools:**Paraver

From timelines to tables



wird cans prome				llayer_mpl_MT_120t.chop1.prv			٢		
	Outside MPI	MPI_Send	MPI_Recv	MPI Isend	MPI Irecv	MPI Waitall	MPI_Bcast	MPI_Reduce	MPL Allr
THREAD 1.113.1	67.6081 %	0.0682 %	9.9182 %	2.5777 %	1.7698-%	5.1676 %	0.5934 %	0.1465 %	-
THREAD 1.114.1	42.8434 %		20.5621 %	1.1947 %	1.0400 %	7.7056 %	1	-	
THREAD 1.115.1	68.6127 %	0.0707 %	9.6223.%	2.2589 %	2.0177 %	S.9825 %	0.5249 %	0.0297 %	
THREAD 1.116.1	74,6039 %	0.0531 %	9.6084 %	2.8813 %	2.5593 %	2.9296 %	0.5095 %	0.0483 %	
THREAD 1.117.1	74,3733 %	0:0691 %	9.7012.%	2.8517 %	2.5240 %	X 🚽 709	MP call profile (i)	Gramacs_bilayes_rapi_M	T_1241.4 (2
THREAD 1.118.1	72,7770 %	0.0545 %	9.5489 %	2.8489 %	2.5353.%	C D 2		H ++    %	
THREAD 1.119.1	66,7994 %	0.0682 %	10.0674 %	2.4206 %	1,9741 %				
THREAD 1.120.1	43.7224 %		20.5273 %	1.1912.96	1.0175.96				
6									
Total	8,012,4546 %	7.3174 %	1.370.5276 %	288,6168 %	253.0137 %	54			
Average	66.7705 %	0.0690 %	11.4211 %	2.4051%	2.1084 %				
Maximum	75.6821 %	0.4390 %	21.2505 %	2.9706 %	2.6369 %				
Minimum	40.5200 %	0.0129 %	8.8583 %	1.1489 %	1.0077 %				
StDev	11.3685 %	0.0474 %	4.0613 %	0.5984 %	0.5406 %				
Avg/Max	0.8822	0.1572	0.5374	0.8095	0.7996				





### **BSC Tools:**Paraver

### One columns per specific value of categorical Control window

<b>C D</b> 30 🔍	🔍 📕 н	H 11 1/2				
	End	MPI_lsend	MPI_Irecv	MPI_Wait	MPI_Alireduce	MPI_Com
THREAD 1.1.1	86,98 %	0,06 %	0,08 %	11,12 %	1,75 %	
THREAD 1.2.1	88,29 %	D,10 %	0,10 %	9,95 %	1,56 %	
THREAD 1.3.1	88,33 %	0,13 %	0,10 %	9,92 %	1,51 %	
THREAD 1.4.1	89,75 %	0,10 %	0,09 %	8,62 %	1,44 %	
THREAD 1.5.1	89,47 %	0,11 %	0,10 %	8,85 %	1,46 %	
THREAD 1.6.1	88,76 %	0,12 %	0,09 %	9,54 %	1,48 %	
THREAD 1.7.1	91,77 %	0,13 %	0,10 %	6,51 %	1,49 %	
THREAD 1.8.1	90,23 %	0,06 %	0,08 %	8,13 %	1,50 %	
THREAD 1.9.1	91,88 %	0,13 %	0,09 %	6,73 %	1,17 %	
HREAD 1.10.1	93,24 %	0,18 %	0,11 %	5,41 %	1,05 %	
HREAD 1.11.1	93,25 %	D,18 %	0,11 %	5,45 %	1,00 %	
THREAD 1.12.1	94,63 %	0,17 %	0,11 %	4,16 %	0,93 %	
THREAD 1.13.1	93,40 %	0,17 %	0,11 %	5,35 %	0,96.%	
THREAD 1.14.1	94,99 %	0,20 %	0,11 %	3,77 %	0,93 %	
HREAD 1.15.1	96,80 %	0,22 %	0,11 %	1,92 %	0,95 %	
THREAD 1.16.1	95,73 %	0,12 %	0,09 %	2,99 %	1,06 %	



Value/color is a statistic computed for the specific thread when control window had the value corresponding to the column

**Relevant statistics:** 

Time, %time, #bursts, Avg. burst time Average of **Data window** 

## MPI calls and profile

- Different types of MPI functions are quantified
- In this case, only the MPI\_Alltoallv and MPI\_Waitany functions represent a significant amount of time with 14.65% and 9.29% respectively.



	Outside MPI	MPI_Recv	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Alltoallv	MPI_Comm_size	MPI_Waitany
Total	41,154.02 %	37.32 %	260.82 %	172.16 %	2,023.25 %	8,438.25 %	163.09 %	5,351.09 %
Average	71.45 %	0.07 %	0.45 %	0.30 %	3.51 %	14.65 %	0.28 %	9.29 %
Maximum	84.51 %	0.18 %	1.15 %	0.81 %	9.70 %	20.07 %	0.31 %	31.12 %
Minimum	51.79 %	0.00 %	0.13 %	0.12 %	0.49 %	8.59 %	0.23 %	1.31 %
StDev	4.91 %	0.04 %	0.18%	0.14 %	1.79 %	3.07 %	0.01 %	4.33 %
Avg/Max	0.85	0.38	0.39	0.37	0.36	0.73	0.93	0.30



Barcelona





## Point-to-point connectivity matrix

- It indicates who communicates with whom
- Almost all point-to-point communications are locally performed between MPI processes neighbours





## **Collective communications**

- Four calls to MPI\_Alltoallv each time step
- The most significant in terms of size and duration is the second one




#### Columns correspond to bins of values of a numeric Control window



Value/color is a statistic computed for the specific thread when control window had the value corresponding to the column

> Relevant statistics: Time, %time, #bursts, Avg. burst time Average of Data window

- Semantic functionality
  - Derived windows
    - Point wise operation
      - $S = \alpha * S^a < op > \beta * S^b$
      - <op> : + , -, \*, /, ...



- Data handling capability
  - Original trace containing all the events
  - Filtering/Burst mode
    - Subset of records in original trace
    - By duration, time, value, event type
    - Trace filtered can be analysed in the same way
    - Also using burst mode from xml file
      - Save only computation bursts longer than a value
  - Cutting
    - All records in a given time interval
    - Only some processes









- Filtering
  - Filter original trace discarding most of the records, only keeping most

relevant information (MPI events can be used for this purpose)

	Cut & Filter	1	= ×	1				
Traces								
Input	ifsMASTER.prv	Browse						
Output	ifsMASTER.filter5.prv	Browse						
	<ul> <li>Load the processed trace</li> </ul>							
	Run application with the processe	ed trace						
Cut/Filter Parame	ters							
Configuration file		Browse						
	Execution chain							
	1 Cutter			5000001				
	✓ 2 Filter	Save		5000001				
	3 Software Counters			50000002	·			
				50000003	MPI events			
Cutter Filter	Software Counters	50100001						
State	Event 🕑 Communication		50100002					
Keep states 50100003								
50100004								
Running								
Not created Unselect all								
Waiting a message								
Bloking Send Min ourst time 0								
The Supplic								
Events								
		Add	$\sim$					
<			>					
			ply					

- Cutting
  - Cut original trace to obtain a fully detailed trace for the time interval considered representative or of interest
  - Use filtered trace to know the area of interest (remember that input must be the original trace)

-	Right click $\rightarrow$ run $\rightarrow$ cutter						
Traces	Cut & Filter 🔹 >						
Input	ifsMASTER.prv Browse	enginal trace					
Output	ITSIVIASTER.chop1.prv Browse						
	<ul> <li>Load the processed trace</li> <li>Run application with the processed trace</li> </ul>						
Cut/Filter Param	eters						
Configuration file	Browse						
	Execution chain  1 Cutter  2 Filter  3 Software Counters						
Cutter Filter Software Counters							
Trace Limits Aied Of Interest							
Cut by time % End 66369058781							
Tasks							
Select Region All Window All Trace							
Trace Options							
< Contract of the second secon							
	<u>○ C</u> ancel  ✓ <u>A</u> pply						

- Configurations for analysis (usr/local/apps/paraver/X.X.X/cfgs)
  - General
    - Including basic views (timelines) and analysis (2D/3D profiles)
  - Counters\_PAPI
    - Hardware counters derived metrics
      - Program: related to algorithm/compilation (instructions, FP ops...)
      - Architecture:related to execution on specific architectures (cache misses...)
      - Performance: metrics reporting rating per time (MIPS, IPC...)
  - MPI  $\rightarrow$  Views and analysis of MPI events
  - OpenMP  $\rightarrow$  Views and analysis of OpenMP events
  - Complete Profile (general\_cfgs)

#### **BSC Tools:Dimemas**

The impossible machine:  $BW = \infty$ , L = 0

- Actually describes/characterizes intrinsic application behavior
  - Load balance problems?



Impact on practical machines?

sendrec

waitall

- Area of study
- Deployment efficiency
- Benchmarking
- Profiling analysis
- Validation







#### • Area of study

- Configuration used (Operational, New algorithms, Global, Parallelization paradigm...)
- Components activated and cyclic patterns
  - IO, ICE, Radiation, MPI, OpenMP
- Area of study
  - 1 complete time step
- Deployment efficiency
- Benchmarking
- Profiling analysis
- Validation







### Types of time step for the practical example

Time steps with radiation are much more expensive due to the extra computation in the grid-point part



ntro Nacional de Supercomputación



#### Structure of a regular time step



- A Inverse transformations
- B Grid-point computations
- C Direct transformations
- **D** Spectral computations



- Area of study (IFS)
  - 24 hours of simulation, T511L137 on CCA (ECMWF)
  - Selected 1 time step: 104 MPI processes + 4 IO (No OpenMP)
  - Metrics collected for large areas of computation automatically









- Area of study (NEMO)
  - 1 day of simulation, ORCA025L91 on MN4 (BSC)
  - Selected the fastest time step automatically
  - 1 time step: 72 MPI processes (No IO, No OpenMP, No SI3)
  - Metrics collected for User functions manually















BSC Barcelona Supercomputing Center Centro Nacional de Supercompulación

- Area of study
- Deployment efficiency
  - Compilation flags
    - Comparing fp options (fast, precise, strict...) and optimization options (OX, vectorization, approximations...)
  - Checking external libraries compilation
  - Debug flags (-g, Optimization reports, -f-instrument-functions...)
- Benchmarking
- Profiling analysis
- Validation







- Area of study
- Deployment efficiency
- Benchmarking
  - Basic Tests to collect Hardware metrics
    - Communications (Latency, Bandwidth, CPU, Parallel Efficiency...)
  - Weak and Strong scaling (MPI, OpenMP, Block processing and Hybrid sets)
  - Comparing optimizations (Double VS Single Precision...)
  - Extrae metrics collection and trace production
- Profiling analysis
- Validation







#### **MPI strong scaling: trace views**



### **Basic Analysis: MPI Strong Scaling**

Computation and parallel efficiency factors for MPI only:

- Good computation scalability and serialization efficiency
- Not very good load balance neither transfer efficiency

Global efficiency

- -- Parallel efficiency
  - -- Load balance
  - -- Communication efficiency
    - -- Serialization efficiency
    - -- Transfer efficiency
- -- Computation scalability
  - -- IPC scalability
  - -- Instruction scalability
  - -- Frequency scalability











### **Basic Analysis: Double P VS Single P**

Overview of the collected raw data:

	108	108		
Runtime (us)	110741508.76	71238767.9		
Runtime (ideal)	105675625.64	68396939.23		
Useful duration (average)	88427932.03	57382830.24		
Useful duration (maximum)	94410288.2	61484222.58		
Useful duration (total)	9196504931.3	5967814345.21		
Useful duration (ideal, max)	94410288.2	61484222.58		
Useful instructions (total)	26798422515714	23201423473963		
Useful cycles (total)	21985000332874	14299301515415		
Overview of the computed mode	l factors: 	108		
Parallel efficiency	79.85%	80.55%		
Load balance	93.66%	93.33%		
Communication efficiency	85.25%	86.31%		
Serialization efficiency	89.34%	89.89%		
Transfer efficiency	95.43%	96.01%		
Computation scalability	100.00%	154.10%		
Global efficiency	79.85%	124.13%		
IPC scalability	100.00%	133.11%		
Instruction scalability	100.00%	115.50%		
Frequency scalability	100.00%   100.23%			
Speedup	1.00	1.55		
Average IPC	1.22	1.62		
Average frequency (GHz)	2.39	2.40		





- Area of study
- Deployment efficiency
- Benchmarking
- Profiling analysis
  - MPI and OpenMP profile summary and Basic Analysis Tool
  - PAPI counters
  - MPI and OpenMP evaluation in detail
  - Clustering and Tracking Tools
  - Sampling and Folding Tools
  - Connection to the code
  - Dimemas Tool
- Validation







### **MPI Profile Summary**

# Parallel and Communication efficiency, Global load balance $\rightarrow$ less than 85%?

Parallel Efficiency

							IFS						
	Outside MPI	MPI_Send	MPI_Recv	MPI_lsend	MPI_Irecv	MPI_Wait	MPI_Barrier	MPI_Alltoallv	MPI_Gatherv	MPI_Comm_rank	MPI_Comm_size	MPI_Bsend	MPI_Waitany
Total	66,578.44 %	1.71 %	773.76 %	646.21 %	239.35 %	12,362.37 %	806.93 %	10,757.31 %	35.56 %	2.49 %	448.23 %	0.81 %	7,746.82 %
Average	66.31 %	0.00 %	0.77 %	0.64 %	0.24 %	12.31 %	0.80 %	10.71 %	0.04 %	0.00 %	0.45 %	0.81 %	7.72 %
Maximum	72.93 %	0.01 %	2.98 %	1.60 %	0.80 %	18.56 %	1.84 %	25.06 %	1.12 %	0.01 %	1.88 %	0.81 %	19.25 %
Minimum	57.05 %	8.00 %	0.01 %	0.08 %	0.07 %	3.11 %	0.00 %	5.25 %	0.00 %	0.00 %	0.16 %	0.81 %	0.31 %
StDev	2.03 %	0.00 %	0.57 %	0.36 %	0.06 %	2.52 %	0.41 %	3.57 %	0.12 %	0.00 %	0.10 %	0 %	3.18 %
Avg/Max	0.91	0.31	0.26	0.40	0.30	0.66	0.44	0.43	0.03	0.34	0.24	1	0.40
MPL Send													· · · · · · · · · · · · · · · · · · ·

Communication Efficiency

**Global Load Balance** 







- PAPI counters collected during the execution
- Some of them are based on other native PAPI counters and derived from the base metrics

	Derived
Instructions	
Cycles	
Useful Duration	Х
Useful Instructions	Х
Useful IPC	Х
Loads	
Stores	
L3/L2/L1_Total_Misses	
L3/L2/L1_MISS_RATIO	Х
FP_OPS	
FP_TOT_INS	
INS_VEC	Х













![](_page_60_Figure_1.jpeg)

#### **MPI evaluation**

![](_page_61_Figure_1.jpeg)

![](_page_61_Picture_2.jpeg)

![](_page_61_Picture_3.jpeg)

![](_page_61_Picture_4.jpeg)

#### **MPI evaluation**

- IPC less than 1 for calculation areas?
- Are there load imbalance regions?

![](_page_62_Figure_3.jpeg)

#### **MPI evaluation**

• Are MPI communications efficient according to the map affinity?

![](_page_63_Figure_2.jpeg)

#### Affinity per node

![](_page_63_Picture_4.jpeg)

Barcelona Supercomputing Center Centro Nacional de Supercomputación

![](_page_63_Picture_6.jpeg)

![](_page_63_Picture_7.jpeg)

### **Clustering Tool**

#### Applying Clustering for an automatic profiling analysis

![](_page_64_Figure_2.jpeg)

- Characterizes computing bursts that are similar and groups them into clusters
- Allows to study the behavior of the clusters separately, identify patterns, etc.

![](_page_64_Picture_5.jpeg)

![](_page_64_Picture_7.jpeg)

## **Tracking Tool**

- A friendly way to quantify and visualize the evolution of the clusters among several traces
- The tool has 2 parts
  - Recognition algorithm of "who-is-who", based on heuristics
  - A visualization GUI
- Examples analyzing multiple traces
  - Scaling number of MPI/OpenMP resources (64 128 256...)
  - Testing different microarchitecture features
  - Changing the problem size
  - Trying different compiler optimizations

![](_page_65_Picture_10.jpeg)

![](_page_65_Picture_12.jpeg)

![](_page_65_Picture_13.jpeg)

### **Tracking Tool**

![](_page_66_Figure_1.jpeg)

![](_page_67_Figure_0.jpeg)

109

108

0.6 0.8 1.0

IPC

PAPI\_TOT\_INS

### **Tracking Tool**

#### Tracking IFS MPI+OMP Strong Scaling

1213

![](_page_67_Figure_3.jpeg)

L1L2

### **Sampling Tool**

- Extrae can be configured to capture performance metrics on a periodic basis using alarm signals and specifying period and variability (10 and 2 respectively for IFS and NEMO tests).
- This means that we will capture samples every 10 ms with a random variability of 2 ms.
- Every sample contains processor performance counters (where every PAPI counter is referred at configured time) and callstack information.

![](_page_68_Picture_4.jpeg)

![](_page_68_Picture_5.jpeg)

![](_page_68_Picture_6.jpeg)

![](_page_69_Figure_0.jpeg)

## **Folding Tool**

- Combine instrumentation and sampling to provide instantaneous performance metrics, source code and memory references. This mechanism receives a trace-file and generates plots showing the fine evolution of the performance.
- The samples collected are gathered from scattered computing regions into a synthetic region by preserving their relative time within their original region so that the sampled information determines how the performance evolves within the region.
- The performance evolution is connected to source code and memory references at the same time.

![](_page_70_Picture_4.jpeg)

![](_page_70_Picture_5.jpeg)

![](_page_70_Picture_6.jpeg)

#### **Folding Tool**

![](_page_71_Figure_1.jpeg)
### **Folding Tool**



### **Folding Tool**



### **Folding Tool**



### **DIMEMAS Tool**

The impossible machine:  $BW = \infty$ , L = 0

- Actually describes/characterizes intrinsic application behavior
  - Load balance problems?



### **DIMEMAS Tool**

Grid point Ideal Network for IFS Semi Grid point calculation Lagrangian execution cmputation (Physics) Actual run MPI call @ fsMASTER.384.1it.prv 843,783 us 1,229,686 us Ideal network call @ D.ideal.ifsMASTER.384.1 t.pr 704,067 us 1,089,969 us Imbalance Transfer sensitive Why does not disappear ?

# **Profiling Methodology**

- Area of study
- Deployment efficiency
- Benchmarking
- Profiling analysis
- Validation
  - Reproducibility Test
  - Validation Test







### Validation

Reproducibility Test: Are your results comparable to the EC-Earth community results?



### The results comparing platforms or configurations:



AMIP platform (Rhino;CCA) comparison Kolmogorov-Smirnov differences of two 5-members ensambles



# Validation Test (NEMO)

- Initial conditions perturbed with white noise in the 3D temperature field.
- Evaluating 53 output variables. ullet

and T estone	erid T heatc.png	rid T mldkz5.png	erid T mldr10 Long	srid T mldr10 ldcv.pm	e erid T saltc png	erid T sbtpng	erid T sopne	erid T sospng	enid T sstdcv.png	erid T taumong	erid T thetao.png
gird_1_cor.pilg	gnu_1_neatc.png	gnu_1_mukz5.png	gild_1_illidi10_1.plig	gnu_1_murro_rucy.pn	g gitu_1_saite.pitg	gita_1_sot.pitg	gild_1_so.plig	gnu_1_sos.png	gitu_1_ssucy.pitg	gnu_1_taumping	gnu_1_metao.png
grid_T_tos.png	grid_T_tosstd.png	grid_T_wfo.png	grid_T_zos.png	grid_T_zosstd.png	grid_U_e3u.png	grid_U_sozohetr.png	grid_U_sozosatr.png	grid_U_tauuo.png	grid_U_uo.png	grid_U_uocetr_eff.png	grid_U_uos.png
grid_U_vozomatrpng	grid_V_e3vpng	grid_V_somehetrpng	grid_V_somesatr.png	grid_V_tauvo.png	grid_V_vo.png	grid_V_vocetr_eff.png	grid_V_vomematr.png	grid_V_vos.png	grid_W_av_ratio.png	grid_W_av_wave.png	grid_W_bflx_iwm.png
							and and and and and and and and				
grid_W_bn2.png	grid_W_difvho.png	grid_W_e3w.png	grid_W_emix_iwm.png	grid_W_pcmap_iwm.png	g grid_W_vovematr.png	grid_W_wo.png	scalar_bgfrchfx.png	scalar_bgfrctem.png	scalar_bgfrcvol.png	scalar_bgheatco.png	scalar_bgheatfx.png
scalar_bgsaline.png	scalar_bgsaltco.png	scalar_bgtemper.png	scalar_bgvole3t.png	scalar_bgvolssh.png	**						
BS	C Supercomp	wing				*					





AND CLIMATE IN EUROPE

# Validation Test (NEMO)

- Initial conditions perturbed with white noise in the 3D temperature field.
- Evaluating 53 output variables.



# Validation Test (NEMO)

### Example: Compiling with -xHost





Barcelona Supercomputing Center Centro Nacional de Supercompulación



AND CLIMATE IN



1				

•		_	
C			
-			
Λ			
4			





BSC Barcelona Supercomputing Center Centro Nacional de Supercomputación



#### Function Timelines





# Border Exchange







- Diagnostic for NEMO:
  - Scalability is constrained by:
    - 1) Algorithms with too much communication
    - 2) Sub-optimal implementation
- Actions taken
  - Improve communication implementation to reduce number of point-to-point messages
  - Reduce number of collectives







· First studies showed that IFS-NEMO coupling was not a big issue



• But it seems that it is when increasing number of cores



BSC has been working successfully with the EC-Earth Technical Working Group to improve the execution of



- A success case: coupling field gathering and OPT option of OASIS coupler for global conservative

#### transformations



- With these optimizations, up to 90% improvement in coupling process can be achieved
- These improvements are now in trunk EC-Earth 3.2.2, substantially benefiting our CMIP6 simulations









BSC has been working successfully with the EC-Earth Technical Working Group to improve the





 A success case: coupling field gathering and OPT option of OASIS coupler for global conservative transformations



- With these optimizations, up to 90% improvement in coupling process can be achieved
- These improvements are now in trunk EC-Earth 3.2.2, substantially benefiting our CMIP6 simulations

- Synchronal point to point communication could be a bottleneck even for only one

message from one master to hundreds of slaves

- Sigcheck method

Date: 1 101	3 A A A A A A A A A A A A A A A A A A A			
THE DO 1 07 1			and the second	
THE R. L 133 1		1 K-4 2	Concerning and the second s	——— <u>}</u> . );
D6740 1.161 1				——— <b>1</b> . 19.
TARTAD 1 101 1				<b>b</b> %
THE R. L				
TAXES 1 157 1				
THEFAD 1. 289. 1				H
D6740 1 321 1				
TAREN 1 152 1				
THERAD 1.385.1				
THEFAD 1.417.1				
THEFAD 1 443 1				and a second
D6740 1 441 1				
THEERO 1.513.1	and the second			
THERE 1.545.1	100 AL 100			
THEEAD 1.577.1				
THEEAD 1.609.1	A		and the second	
THE 240 1.641.1	1 2 2 E			
THEE20 1.073.1				
THREAD 1.705.1	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1			
THEEAD 1.737.1	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1			
THREAD 1.769.1				
THREAD 1.000.1				
THEERO 1.633.1				
THREAD 1.865.1				
THEEAD 1.897.1				
THEEAD 1.929.1				
THREAD 1.961.1				
THEE20 1.992.1	11 11		and the second	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1

MPI call protile @ itsMASTER.prv #8

	89 (sigcheck.F90, ifsMASTER)	95 (sigcheck.F90, ifsMASTER)	111 (sigcheck.F90, ifsMASTER)
5	0.81 %	442.74 %	806.93 %
5	0.81 %	0.44 %	0.80 %
5	0.81 %	1.73 %	1.84 %
5	0.81 %	0.00 %	0.00 %
	0 %	0.45 %	0.41 %
ī	1	0.25	0.44

- Using one asynchronal collective communication this time is

### reduced almost to 0







### • Hybrid Test (128 MPI+4 OpenMP, Total: 512)

**OpenMP Parallel Regions** 









Small OpenMP parallel Regions



128 MPI processes and 4 OpenMP threads per process

128 MPI processes and 4 OpenMP threads per process













Barcelona Supercomputing Center Centro Nacional de Supercomputación



**ESIVACE** CENTRE OF EXCELLENCE IN SIMULATION OF WEATHER AND CLIMATE IN FUROPE

# Thank you

The research leading to these results has received funding from the EU H2020 Framework Programme under grant agreement H2020 GA 675191.

The content of this presentation reflects only the author's view. The European Commission is not responsible for any use that may be made of the information it contains.

### mario.acosta@bsc.es