

Earth System Modelling on System-level Heterogeneous Architectures: EMAC (version 2.42) on the Dynamical Exascale Entry Platform (DEEP)

M. Christou¹, T. Christoudias¹, J. Morillo², D. A. Mallon³ and H. Merx^{1, 4}

[1]{The Cyprus Institute, Nicosia, Cyprus}

[2]{Barcelona Supercomputing Center, Barcelona, Spain}

[3]{Jülich Supercomputing Centre, Jülich, Germany}

[4]{Max Planck Institute for Chemistry, Mainz, Germany}

Correspondence to: T. Christoudias (christoudias@cyi.ac.cy)

Abstract

We examine an alternative approach to heterogeneous cluster-computing in the many-core era for Earth System models, using the European Centre for Medium-Range Weather Forecasts Hamburg (ECHAM)/Modular Earth Submodel System (MESSy) Atmospheric Chemistry (EMAC) model as a pilot application on the Dynamical Exascale Entry Platform (DEEP). A set of autonomous coprocessors interconnected together, called Booster, complements a conventional HPC Cluster and increases its compute performance, offering extra flexibility to expose multiple levels of parallelism and achieve better scalability. The EMAC model atmospheric chemistry code (Module Efficiently Calculating the Chemistry of the Atmosphere (MECCA)) was taskified with an offload mechanism implemented using OmpSs directives. The model was ported to the MareNostrum 3 supercomputer to allow testing with Intel Xeon Phi accelerators on a production-size machine. The changes proposed in this paper are expected to contribute to the eventual adoption of Cluster-Booster division and Many Integrated Core (MIC) accelerated architectures in presently available implementations of Earth System Models, towards exploiting the potential of a fully Exascale-capable platform.

1 **1 Introduction**

2 The ECHAM/MESSy Atmospheric Chemistry (EMAC) model is a numerical chemistry and
3 climate simulation system that includes sub-models describing tropospheric and middle
4 atmosphere processes and their interaction with oceans, land and human influences (Jöckel et
5 al., 2010). It uses the second version of the Modular Earth Submodel System (MESSy2) to link
6 multi-institutional computer codes. The core atmospheric model is the 5th generation European
7 Centre for Medium Range Weather Forecasts Hamburg general circulation model (ECHAM5,
8 Roeckner et al., 2003, 2006).

9 The EMAC model runs on several platforms, but it is currently unsuitable for massively parallel
10 computers, due to its scalability limitations and large memory requirements per core. EMAC
11 employs complex Earth-system simulations, coupling a global circulation model (GCM) with
12 local physical and chemical models. The global dynamical processes are strongly coupled and
13 have high communication demands while the local physical processes are inherently
14 independent with high computation demands. This heterogeneity between different parts of the
15 EMAC model poses a major challenge when running on homogeneous parallel supercomputers.

16 We test a new approach for a novel supercomputing architecture as proposed by the DEEP
17 project (Eicker et al., 2013, 2015, Mallon et al., 2012, 2013, Suarez et al, 2011), an innovative
18 European response to the Exascale challenge. Instead of adding accelerator cards to Cluster
19 nodes, the DEEP project proposes to use a set of interconnected coprocessors working
20 autonomously (called Booster), which complements a standard Cluster. Together with a
21 software stack focused on meeting Exascale requirements—comprising adapted programming
22 models, libraries and performance tools—the DEEP architecture enables unprecedented
23 scalability. The system-level heterogeneity of DEEP, as opposed to the common node-level
24 heterogeneity, allows users to run applications with kernels of high scalability alongside kernels
25 of low scalability concurrently on different sides of the system, avoiding at the same time over
26 and under-subscription.

27 The Cluster–Booster architecture is naturally suited to global atmospheric circulation–
28 chemistry models, with global components running on the Cluster nodes exploiting the high-
29 speed Xeon processors and local components running on the highly-parallel Xeon Phi co-
30 processors. By balancing communication versus computation the DEEP concept provides a new
31 degree of freedom allowing us to distribute the different components at their optimal
32 parallelisation.

1 **2 Overview of application structure**

2 The EMAC model comprises two parts, the dynamical base model ECHAM, using a nonlocal,
3 spectral algorithm with low scalability, and the modular framework MESSy, linking local
4 physical and chemical processes to the base model, with high scalability. While the number of
5 processors used for the base model is limited by the non-local spectral representation of global
6 dynamical processes, local physical and chemical processes described by framework submodels
7 run independently from their neighbours and present very high scalability.

8 **2.1 Phases**

9 The ECHAM base model runs in parallel in the distributed-memory paradigm using the
10 Message Passing Interface (MPI, Aoyama et al., 1999) library for communication; the MESSy
11 framework inherits the parallel decomposition defined by the base model. While ECHAM has
12 been shown to be able to exploit the shared-memory paradigm using the Open Multi-Processing
13 (OpenMP) library (Dagum et al., 1998), no such effort had been undertaken for the MESSy
14 model so far.

15 It is, however, currently not possible to delegate the whole MESSy subsystem to full multi-
16 threaded execution as some physical processes are naturally modelled in a column-based
17 approach, and are strongly dependent on the system states at their vertically adjacent grid points.
18 The implementation of submodels simulating these processes consequently relies on the column
19 structure inherited from the base model, e.g. sunlight intensity at lower grid points depending
20 on the absorption at higher grid points, and precipitation depending on the flux of moisture from
21 vertically adjacent grid cells

22 Describing gas phase chemical kinetics, the MESSy submodel MECCA executes independently
23 of its physical neighbours and is not limited by vertical adjacency relations. As more than half
24 of the total run-time is spent in MECCA for a typical model scenario, it seems adequate to
25 concentrate on the MECCA kernel with strong algorithmic locality and small communication
26 volume per task. As sketched in Figure 1 the current implementation of MECCA, developed in
27 the DEEP project, is delegated to the Booster using a task-based approach while both ECHAM
28 and the remaining MESSy submodels are executed on the Cluster in the distributed-memory
29 paradigm.

1 **2.2 Scalability dominant factors**

2 Implementing a spectral model of the dynamical state of the atmosphere, the ECHAM phase
3 comprises six transform and six transposition operations in each time step. The data in memory
4 for each time step (data size scales with the square of the model horizontal resolution) is
5 transposed in an all-to-all communication pattern, and this phase is dominated by network
6 bandwidth.

7 Figure 2 displays one model cycle traced with Extrae/Paraver (Extrae 2015, Paraver 2015)
8 starting with the end of the grid point calculations of the last time step calculated—in which
9 most processors are already idle (orange) due to load imbalance and waiting for process 14
10 (blue) to finish running. This is followed by the transpositions in the beginning of a new time
11 step, and Fourier and Legendre transformations (magenta), which execute simultaneously as
12 further analysis showed. After the transpositions a short interval with all processors running
13 (blue) can be identified with the time step integration in spectral space, followed by the inverse
14 transformations and transpositions and transport calculations in ECHAM.

15 While the pattern described so far repeats towards the end of the displayed interval, the major
16 fraction of the time step is spent without communication, running (blue) or waiting (orange) in
17 calculations in MESSy in grid space. The MESSy phase comprises some 30 submodels that are
18 tightly coupled by exchanging the atmospheric observables using global variables.
19 Investigations during the first phase of the project determined the load imbalance visible in
20 Figure 2 to be caused by chemical processes computed in the MECCA submodel.

21 The observed load imbalance is one of the main factors determining application scalability. It
22 is caused by an adaptive time-step integrator solving a system of differential equations. As the
23 stiffness of these equations representing photochemical reactions varies by up to two orders of
24 magnitude due to changes in the intensity of sunlight, the adaptive integrator demands varying
25 amounts of run time accordingly.

26 In the MECCA phase the algorithmically complex adaptive time-step differential equation
27 integrator operates on chemical concentrations of a total data size of the order of a few kilobytes
28 per grid point. Yet, as observed using Scalasca (Scalasca 2015), this phase consumes the major
29 proportion of the total execution time (76%), it is compute-bound and an obvious candidate for
30 offloading to accelerators. It should be noted though, that in a regular architecture, accelerating
31 this highly parallel phase will not eliminate the load imbalance.

1 To test the model scalability out-of-the-box, the EMAC application has been ported to the
2 JUDGE cluster at the Jülich Supercomputing Centre (JSC), and a representative benchmark
3 with a horizontal resolution of 128 grid points in longitudinal and 64 grid points in latitudinal
4 direction with 90 vertical levels and a spin-up period of 8 simulated months has been compiled,
5 frozen and packaged to be used for measurements. Table 1 details the experimental setup for
6 the results shown in this section. The default E5M2 KPP chemistry batch option, along with
7 model namelist setup (NML_SETUP) E5M2/02b_qctm, without dynamics nudging and with
8 the diagnostic submodels D14CO, DRADON, S4D and VISO switched off, is used throughout
9 this work. All disk output channels were turned off for timing purposes.

10 The EMAC strong scaling was benchmarked with different numbers of processors on JUDGE
11 in order to determine the run time behaviour of the total application. As shown in Figure 3 the
12 application scales up to 384 processes (16 nodes x 24 MPI processes each), at higher numbers
13 the performance decreases. Parallel execution speed is determined by the balance of three
14 factors: computation, communication, and load imbalance. The benchmarking setup for the
15 JUDGE cluster can be seen in Table 2. While the computational resources increase with
16 additional processors and therefore increase the application performance, communication
17 demands diminish the positive effect of the additional processors. Additionally, increasing the
18 granularity of the total workload also increases the load imbalance.

19 While the number of processors used for the distributed-memory part of the code is limited by
20 the scalability of the non-local representation of global dynamical processes in ECHAM, the
21 local processes in MESSy running independently from their neighbours scale very well. The
22 MESSy subsystem has not been designed for multi-threaded execution, though, and contains
23 non-local code due to characteristics of the physical processes and algorithmic design decisions.
24 Some physical processes are naturally modelled in a column-based approach, because they are
25 strongly dependent on the system states at vertically adjacent grid points, e.g. sunlight intensity
26 at lower grid points depending on the absorption at higher grid points, and precipitation
27 depending on the flux of moisture from vertically adjacent grid cells. Sub-models simulating
28 these processes consequently rely on the column structure implemented in the current model.

29 In the existing distributed-memory parallel decomposition, the three-dimensional model grid is
30 split horizontally using two run-time parameters, setting the number of processes in latitudinal
31 and longitudinal direction. As work is distributed independently for each direction, a
32 rectangular decomposition is obtained. In ECHAM5, two such rectangular sets of grid-points,

1 symmetric with respect to the equator and balancing the load distribution between the Earth's
2 hemispheres, are assigned to one processor. Further, the first transposition in the dynamical
3 core reorders the variables assigned to a processor to rectangular stripes extending along the
4 full latitudes. Assigning the two run-time parameters mentioned above so that the rectangular
5 sets are elongated in latitudinal direction reduces the inter-processor communication in the first
6 transposition.

7 The physical load-imbalance, caused by photo-chemical processes in the lower stratosphere and
8 natural and anthropogenic emissions (lightning, soil bacteria, fuel combustion—motor vehicles,
9 industrial, and utility), appears in the run time spent for each grid point when examining the
10 benchmark calculations. In Figure 4 the maximal MECCA kernel execution wall-time for one
11 grid point in each column differs by up to a factor of four. The load imbalance is caused by the
12 adaptive time-step integrator solving the differential equations that describe the chemical
13 equations computed in the MECCA submodel.

14 At high levels of parallelisation, the load imbalance becomes a limiting factor, and the factors
15 determining scalability in absolute numbers in Figure 5 are both communication and
16 computation. For the ECHAM phase (blue), when scaling to beyond 8 nodes the
17 communication demands of the underlying spectral model involving several all-to-all
18 communication patterns start to dominate. Computing time for MESSy is still decreasing and
19 the computing time for MECCA decreases stronger than MESSy (note logarithmic scale and
20 distinction of MECCA from MESSy). For reference, the MareNostrum 3 heterogeneous
21 compute nodes have each 2x8 CPU cores (MPI) and 2x Xeon Phi accelerators.

22 In Figure 6 the point at which communication and computation require equal times around 8
23 nodes is clearly apparent; this ultimately limits the maximum number of cores that can be used
24 with high efficiency and 16 nodes (equivalent to 256 cores) is commonly used in production
25 runs of the EMAC atmospheric model as a scientific application to balance efficiency and total
26 required wall time.

27 **3 Model Developments**

28 **3.1 Intranode taskification**

29 The EMAC model atmospheric chemistry code (MECCA) was taskified using OmpSs ([Bueno](#)
30 [et al., 2011, 2012](#), [Duran et al., 2011](#), [Sainz et al., 2014](#)) directives. OmpSs allows the user to
31 specify inputs and outputs for blocks of code or functions, giving enough information to the

1 Nanos++ runtime system to construct a dependency graph. This dependency graph reflects at
2 all moments which tasks are ready to be executed concurrently, and therefore the programmer
3 does not have to explicitly manage the parallelisation. The concept of tasks and task
4 dependencies has also been adopted in the OpenMP 4.0 standard (OPENMP 4.0, 2013). Since
5 in MECCA each gridpoint is computed independently of its neighbours, this part of the code is
6 in principle embarrassingly parallel, with no communication or inter-task dependencies
7 involved. The MECCA submodel was refactored through the creation of computational kernels
8 for intranode parallelisation with shared-memory tasks. All changes are in the automatically
9 generated MESSy MECCA KPP source code and no additional changes are needed in the
10 MESSy submodel interface layer (SMIL) or in the core layer (SMCL). Since the source is
11 automatically generated, the changes have to be propagated to the generator after the KP4
12 mechanism runs (applied in order to remove the indirect indexing and expand the original KPP
13 code by an additional dimension to enable a better performance due to better cache usage). The
14 propagation is still work in progress and falls outside the scope of this manuscript.

15 The new integration subroutine performs identical functions as the `kpp_integrate` subroutine
16 but on unique elements:

```
17 !$OMP TARGET DEVICE(SMP) COPY_IN(k, temp, press, cair, khet_st,  
18 khet_tr, jx, time_step_len, atol, rtol, icntrl, rcntrl)  
19 COPY_INOUT(conc)
```

```
20 SUBROUTINE kpp_integrate_kernel (k, conc, temp, press, cair,  
21 khet_st, khet_tr, jx, dt, atol, rtol, icntrl, rcntrl)
```

22 The statements `COPY_INOUT` and `COPY_IN` refer to the OmpSs specification and declare
23 memory dependence and declare task dependencies. These dependencies are used by the
24 Nanos++ runtime system to construct a dependency graph which ensures that the tasks are
25 executed concurrently when their dependencies are met and there are free threads. The OmpSs
26 pragma targets symmetric multiprocessing (**SMP**), symmetric multiprocessor system hardware
27 and software architecture where each identical processor unit connects to a single, shared main
28 memory and has full access to all I/O devices.

29 The new version of EMAC, running ECHAM with MPI processes and MECCA with shared-
30 memory OmpSs tasks outperforms the old EMAC using pure MPI, and continues to scale

1 beyond the region where the original implementation scaling performance plateaus. This can
2 be seen in Figure 4, which shows the performance using multi-threading on the DEEP Cluster.

3 **3.2 Internode taskification**

4 In DEEP, OmpSs has been extended to support offloading tasks to remote nodes (Beltran et al.,
5 2015). This mimics the behaviour of other accelerator APIs that move data from the host to the
6 device, compute in the device, and return the results to the host. However, OmpSs adds two
7 very important features: i) it allows offloading to remote nodes, not just locally available
8 coprocessors/accelerators, which is a key functionality to effectively use the Booster; and ii) it
9 allows using the Booster as a pool of coprocessors, so tasks can be offloaded to any Booster
10 node with enough free cores. The latter enables to eliminate the load-imbalance caused by
11 sunlight gradients in MECCA. The second source of imbalance by heterogeneous reactions is
12 also automatically alleviated by the dynamical load balance using the massive parallelisation in
13 the Booster. Our proposed solution is agnostic to the specific origin of load imbalance in the
14 chemistry calculation.

15 In a shared-memory taskification the data is already shared between threads, and no memory
16 copies are necessary. However, in DEEP, to leverage the Booster, this data has to be copied to
17 the Booster nodes. Keeping that in mind, the new task-based MECCA implementation was
18 optimised and the memory and network footprint of the distributed-memory offloading was
19 reduced by three orders of magnitude. To minimise the memory footprint for offloaded tasks,
20 the number of computational grid elements issued to MESSy is further split into individual
21 elements for each task, by rearranging the grid point arrays in each time step to implement data
22 locality at the grid-point level, resulting in a reduction of the total memory footprint from 2.7
23 MB down to 6.3 KB for each task. This was the result of refactoring both the data and code
24 structures in MECCA.

25 In order to create multiple threads integrating the chemical equations simultaneously, all data
26 describing the state of a single grid point have been identified and separated into local variables.
27 Using these variables an integrator kernel has been created that can be offloaded onto threads
28 running on the main processor or hardware accelerators. The integrator kernel needs the
29 variables describing the local state of the atmosphere and the integrator parameters defining the
30 solution of the chemical equations. As the chemical mechanism is compiled by the Kinetic Pre-
31 Processor (KPP) some module variables contain the data of a set of grid points that is accessed

1 using the index K inside the integrator. In order to maintain the code for the integration of
2 different chemical mechanisms created by the KPP compiler these input variables are passed in
3 their original size together with the column grid point index K:

```
4 !OMP TASK FIRSTPRIVATE (k) INOUT (conc (k, :) ) IN (temp, press,  
5 cair, khet_st, khet_tr, jx, time_step_len, atol, rtol, icntrl,  
6 rcntrl) ONTO ( )
```

7 The keywords INOUT and IN refer to the OmpSs specification and declare task dependencies.
8 ONTO sets the offload target—left blank it targets any available accelerator device. The input
9 arrays have been reduced in the dimension of NBL (see below), and now carry the memory of
10 a single grid point. The only two-dimensional array left is conc, and only one slice (k,:) is passed
11 to and from containing the concentrations of all chemical species. Prior to passing each element
12 all two-dimensional arrays above are transposed in the outermost dimension, to have sequential
13 memory regions, as in the Fortran language specification. The compiler directives are
14 implemented as pragmas (comments in the source code), automatically controlled with
15 definitions at compile time (being invoked only if SMP/Mercurium is present), thus no
16 additional user input is required. The kernel data sizes depend on parameters in the EMAC
17 environment $NBL = KPROMA * NLEV$. The block length NBL corresponds to the vector length
18 used in the ECHAM base model for grid-point calculations. A block in ECHAM comprises all
19 grid points in a slice of the model grid extending along a virtual longitude of length KPROMA
20 and all model levels NLEV. KPROMA is a run-time parameter and NLEV corresponds to the
21 model vertical resolution. The vector length in the integrator kernel is defined by the user based
22 on the capabilities of the hardware at compile time. The total size of data available to a single
23 shared-memory task depends on the KPP control parameters and number of species (NSPEC).
24 As an example, using the benchmark parameters for the chemical mechanism each task has
25 access to 2.7 Mbytes of shared main memory.

26 When offloading tasks, to maximise massive parallelisation, the parameter NBL is reduced at
27 runtime to $NBL = 1$, by redefining the grid point arrays, resulting in a total memory footprint
28 of just 3588 bytes for each task. At the benchmark resolution of T42L90MA a total number of
29 737 280 independent tasks is generated in each time step resulting in a total data size of 2.5
30 Gbytes. Creating the tasks requires all this memory to be transferred from the Cluster to the
31 Booster.

1 At a typical parallelisation of 256 processes for the less-scaling spectral parts of the model each
2 MPI process creates $NBL = 2880$ tasks and transfer 9.9 Mbytes per process in each simulation
3 time step. The data returned from the OmpSs tasks is smaller in size as it only comprises the
4 concentrations of each chemical species in the CONC array. For the benchmark chemical
5 mechanism with $NSPEC = 139$ this amounts to 1112 bytes per task. The correspondig total
6 amount of data per time step is 782 Mbytes, with each MPI process receiving 3.0 Mbytes.

7 Additionally, the distributed-memory offloading code was redesigned to exploit shared memory
8 within the Xeon Phi many-core processors by nesting an OmpSs shared-memory region within
9 Cluster-to-Booster tasks encompassing variable, runtime-defined number of individual
10 gridpoint calculations. The runtime parameter can be tweaked to change the total memory
11 required by the number of gridpoints per task to fit within the device shared memory, to
12 minimise runtime overhead (avoiding many tasks spawned over the network, which has a huge
13 penalty). Thus, the number of tasks to be sent to the Booster can be controlled and optimised
14 for each architecture, and host-specific configuration allows for optimum task size based on
15 bandwidth, reducing task communication overheads.

16 With this approach, the specifics of the DEEP system architecture, and in particular the
17 hardware present in MIC coprocessors is exploited by massively parallelising the chemistry
18 calculations at the gridpoint level and offloading to the Booster exposing a significant amount
19 of thread parallelism. At the same time the load imbalance observed in MECCA is
20 automatically alleviated through OmpSs' dynamic load balancing by selecting a sufficiently
21 fine task size and decoupling the model-domain location of the grid point from the task
22 execution on the physical CPU.

23 **3.3 Kernel refactoring**

24 The computational kernel of MECCA conceptually accesses several single elements of the
25 global arrays describing the state of the atmosphere. In order to create data locality, in a first
26 step the module variables used in the current version of MECCA to represent the model state
27 were refactored within the `mecca_physc` function to be passed explicitly to the `kpp_integrator`
28 routine that computes the chemical equations of one vector of grid points. Subsequently, to
29 increase memory adjacency, some fields were transposed from the memory order optimised for
30 vectorisation along longitudes as given by the ECHAM core model to arrays with adjacent
31 elements representing chemical concentrations and reaction rate constants. In a third step, the

```

1  computation of the equations of a single grid point were encapsulated into a new
2  kpp_integrator_kernel function that calls the functions update_rconst and integrate, which are
3  created automatically by the KPP compiler. This innermost routine is called in a loop over all
4  grid points in a vector along the virtual longitude defined by the ECHAM infrastructure,
5  DO k = 1, vl_glo
6      CALL kpp_integrate_kernel (conc (:, k), temp (k), press (k),
7  cair (k), khet_st (:, k), khet_tr (:, k), jx (:, k), time_step_len,
8  atol, rtol, icntrl, rcntrl, wtime (k))
9  END DO

```

10 By maintaining the interfaces both to the MESSy infrastructure and the code implementing the
11 KPP integrator, the changes to the EMAC code base were kept minimal, and a transfer of the
12 changes to other MESSy submodels using the KPP integrator should be feasible.

13 **4 Attainable Performance**

14 At the time of writing this manuscript the DEEP Booster is in the bring-up phase, and not
15 available to users. In order to project the performance of the full DEEP System, Xeon-based
16 measurements on the DEEP Cluster were combined with Xeon Phi-based measurements on
17 MareNostrum 3. The “Pure MPI” nominal time on the DEEP Cluster for all phases (red line),
18 the time for each phase, and the theoretical performance when offloading to the Booster are
19 shown in Figure 5. The DEEP Cluster reference data weighted by the relative factors for each
20 phase derived from the metrics measurements exhibit a performance maximum for the base
21 model (ECHAM) and MESSy (excluding MECCA) at 8 nodes (light blue line), representing a
22 good estimate for the optimal parallelisation of that phase on the Cluster. This estimate of 375
23 s per simulated day for the low-scaling Cluster phases was used to extrapolate the attainable
24 performance (green line), merging this result at 8 nodes, with the Xeon Phi data retrieved from
25 MareNostrum 3, where benchmarks using one node had been run with varying numbers of
26 processing elements within one Xeon Phi processor.

27 While the number of Booster nodes required to attain similar performance to the original
28 distributed-memory based implementation corresponds to regular accelerator architectures with
29 individual boosters directly attached to cluster nodes, the projected DEEP performance (green
30 line) scales beyond the optimal performance achieved so far. The EMAC atmospheric
31 chemistry global climate model seems therefore well suited to exploit an architecture providing

1 considerable more hardware acceleration than provided by regular systems. To find the
2 attainable performance time $T_{\text{attainable}}$ when running the ECHAM+MESSy on the Cluster and
3 MECCA on the Booster, varying the number of Booster nodes, the following equation is used:

$$4 \quad T_{\text{attainable}} = T_{\text{ECHAM+MESSy}}^{\text{min}} + T_{\text{MECCA}}^{\text{Booster}}$$

5 The projected attainable performance that outperforms the pure-MPI conventional cluster
6 paradigm at higher core count (depicted here as the number of Booster nodes, while keeping
7 the ECHAM/MESSy MPI part on 8 Cluster nodes for optimal performance $T_{\text{ECHAM+MESSy}}^{\text{min}}$) is
8 also shown in Figure 5.

9 **5 Conclusions**

10 The global climate model ECHAM/MESSy Atmospheric Chemistry (EMAC) is used to study
11 climate change and air quality scenarios. The EMAC model is constituted by a nonlocal
12 dynamical part (ECHAM) with low scalability, and local physical (MESSy) and chemical
13 (MECCA) processes with high scalability. The model's structure naturally suits the DEEP
14 Architecture using the Cluster nodes for the nonlocal part and the Booster nodes for the local
15 chemistry processes. Different implementations of the code's memory and workload divisions
16 were developed and benchmarked to test different aspects of the achievable performance on the
17 proposed architecture. The use of the OmpSs API largely frees the programmers of
18 implementing the offloading logic and, given that EMAC is developed and used in a large
19 community working on all aspects of the model, can facilitate adoption of the concept in the
20 MESSy community.

21 The chemistry mechanism was taskified at the individual gridpoint level using OmpSs
22 directives. The chemistry code was refactored to allow for memory adjacency of vector
23 elements. The OmpSs taskification with remote offload allows for massive task parallelisation
24 and the implementation of optional two-stage offload to control Cluster-Booster task memory
25 size and optimum bandwidth utilisation.

26 The computational load imbalance arising from a photochemical imbalance is alleviated at
27 moderate parallelisation by assigning grid points with differing run times to each process and
28 distributing the load over all processes. Due to the physical distribution of sunlight this load
29 balancing does not require an explicit algorithm at moderate parallelisation; instead, the implicit
30 assignment of the model grid in rectangular blocks suffices for this purpose. At higher numbers
31 of processors this implicit load-balancing decreases and the resulting load imbalance has to be

1 solved by active balancing. The dynamic scheduling provided by the OmpSs run-time system
2 balances the computational load without a possible, but expensive prediction for the current
3 time step.

4 With these approaches, the specifics of the DEEP System architecture, and in particular the
5 hardware present in MIC coprocessors can be exploited by massively parallelising the
6 chemistry calculations at the gridpoint level and offloading to the Booster exposing a significant
7 amount of thread parallelism. At the same time the load imbalance observed in MECCA will
8 be automatically alleviated through dynamic load balancing by minimising the individual task
9 size to one grid box and decoupling the model-domain location from the task execution on the
10 physical CPU, and transferring it to any available core on the Booster.

11 Benchmark projections based on available hardware running the DEEP software stack suggest
12 that the EMAC model requires the large numbers of Xeon Phi accelerators available in the
13 DEEP architecture to scale beyond the current optimal performance point and exploit Amdahl's
14 law with the highly scalable gridpoint calculations while capitalising on the high performance
15 and fast communication for the spectral base model on Intel Xeon processors.

16 In the longer term, when the optimal speedup of the chemical kinetics computation is
17 approached, the overall parallel speedup of EMAC will be limited by the remaining non-
18 accelerated submodels. To achieve further speedup, additional computationally expensive
19 submodels like the aerosol submode GMXe could be addressed. GMXe is column-bound and
20 calls KPP up to four times during one time step (for grid/subgrid scale liquid/ice clouds) with
21 load imbalance caused by the distribution of clouds over the model domain. In principle, the
22 MECCA implementation is directly applicable to the case of SCAV. The actual performance
23 gain would be dependent on the exact setup and remains to be tested.

24 The changes proposed in this paper are expected to contribute to the eventual adoption of MIC
25 accelerated architectures for production runs, in presently available implementations of Earth
26 System Models, towards exploiting the potential of a fully Exascale-capable platform.

27 **Code Availability**

28 The Modular Earth Submodel System (MESSy) is continuously further developed and applied
29 by a consortium of institutions. The usage of MESSy and access to the source code is licensed
30 to all affiliates of institutions which are members of the MESSy Consortium. Institutions can
31 become a member of the MESSy Consortium by signing the MESSy Memorandum of

1 Understanding. More information can be found on the MESSy Consortium Website
2 (<http://www.messy-interface.org>). Changes to the source code to implement system-level
3 heterogeneous offloading are currently specific to the Mercurium compiler, Nanos++ runtime
4 and proprietary Parastation MPI and are hosted on the DEEP project version control repository.

5

6 **Acknowledgements**

7 The research leading to these results has received funding from the European Community's
8 Seventh Framework Programme (FP7/2007-2013) under Grant Agreement n° 287530.

9

1 **References**

- 2 Aoyama, Y. and Nakano, J., *RS/6000 SP: Practical MPI Programming*, 1999.
- 3 Beltran, V., Labarta, J. and Sainz, F., *Collective Offload for Heterogeneous Clusters*, IEEE
4 International High Performance Computing (HiPC), Bangalore, India, 2015.
- 5 Bueno, J., Planas, J., Duran, A., Badia, R.M., Martorell, X., Ayguade, E. and Labarta, J.,
6 *Productive Programming of GPU Clusters with OmpSs*, Parallel & Distributed Processing
7 Symposium (IPDPS), 2012 IEEE 26th International, 557–568, 2012.
- 8 Bueno J., Martinell L., Duran A., Farreras M., Martorell X., Badia R. M., Ayguade E. and
9 Labarta J., *Productive cluster programming with OmpSs*, Euro-Par 2011 Parallel Processing,
10 Lecture Notes in Computer Science **6852**, 555–566, 2011.
- 11 Dagum, L. and Menon, R.: OpenMP: an industry standard API for shared-memory
12 programming, Computational Science & Engineering, IEEE **5**, 1, 46–55, 1998.
- 13 Damian, V., Sandu, A., Damian, M., Potra, F. and Carmichael, G.R.: *The Kinetic PreProcessor*
14 *KPP—A Software Environment for Solving Chemical Kinetics*, Computers and Chemical
15 Engineering **26**, 11, 1567–1579, 2002.
- 16 Duran, A., Ayguadé, E., Badia, R., M., Labarta, J., Martinell, L., Martorell, X., and Planas, J.,
17 *OmpSs: a proposal for programming heterogeneous multi-core architectures*, Parallel
18 Processing Letters **21**, No. 02, 173–193, 2011.
- 19 Eicker, N., Lippert, T., Moschny, T. and Suarez, E., *The DEEP project: Pursuing cluster-*
20 *computing in the many-core era*, Proc. of the 42nd International Conference on Parallel
21 Processing Workshops (ICPPW) 2013, Workshop on Heterogeneous and Unconventional
22 Cluster Architectures and Applications (HUCAA), Lyon, France, 885–892, 2013.
- 23 Eicker, N., Lippert, T., Moschny, T. and Suarez, E., *The DEEP Project—An alternative*
24 *approach to heterogeneous cluster-computing in the many-core era*, Concurrency and
25 Computation, 2015.
- 26 Extrae, Barcelona Supercomputing Center, <https://www.bsc.es/computer-sciences/extrae>, Last
27 Access: 23/11/2015.

- 1 Jöckel, P., Kerkweg, A., Pozzer, A., Sander, R., Tost, H., Riede, H., Baumgaertner, A., Gromov
2 S. and Kern, B., *Development cycle 2 of the Modular Earth Submodel System (MESSy2)*,
3 *Geoscientific Model Development* **3**, 717–752, 2010.
- 4 Mallon, A. D., Lippert, T., Beltran, V., Affinito, F., Jaure, S., Merx, H., Labarta, J.,
5 Staffelbach, G., Suarez, E. and Eicker, N., *Programming Model and Application Porting to the*
6 *Dynamical Exascale Entry Platform (DEEP)*, Proceedings of the Exascale Applications and
7 Software Conference, Edinburgh, Scotland, UK., 2013.
- 8 Mallon, A., D., Eicker, N., Innocenti, M.E., Lapenta, G., Lippert, T. and Suarez, E., *On the*
9 *Scalability of the Cluster-Booster Concept: a Critical Assessment of the DEEP Architecture*,
10 FutureHPC '12, Proceedings of the Future HPC Systems: the Challenges of Power-
11 Constrained Performance, ACM New York, 2012.
- 12 OpenMP Application Program Interface Version 4.0 – July 2013, OpenMP Architecture
13 Review Board, <http://www.openmp.org/mp-documents/OpenMP4.0.0.pdf>, Last Access:
14 23/11/2015.
- 15 Paraver, Barcelona Supercomputing Center, [https://www.bsc.es/computer-](https://www.bsc.es/computer-sciences/performance-tools/paraver)
16 [sciences/performance-tools/paraver](https://www.bsc.es/computer-sciences/performance-tools/paraver), Last Access: 23/11/2015.
- 17 Roeckner, E., Brokopf, R., Esch, M., Giorgetta, M., Hagemann, S., Kornblueh, L., Manzini,
18 E., Schlese, U. and Schulzweida, U., *Sensitivity of simulated climate to horizontal and*
19 *vertical resolution in the ECHAM5 atmosphere model*, *J. Climate* **19**, 3771–3791, 2006.
- 20 Roeckner, E., Bäuml, G., Bonaventura, L., Brokopf, R., Esch, M., Giorgetta, M., Hagemann,
21 S., Kirchner, I., Kornblueh, L., Manzini, E., Rhodin, A., Schlese, U., Schulzweida, U. and
22 Tompkins, A., *The atmospheric general circulation model ECHAM 5. PART I: Model*
23 *description*, *Report/MPI für Meteorologie* **349**, 2003.
- 24 Sainz, F., Mateo, S., Beltran, V., Bosque, J., L., Martorell, X. and Ayguadé, E., *Leveraging*
25 *OmpSs to Exploit Hardware Accelerators*, International Symposium on Computer
26 Architecture and High Performance Computing, 112–119., 2014.
- 27 Sander, R., Baumgaertner, A., Gromov, S., Harder, H., Jöckel, P., Kerkweg, A., Kubistin, D.,
28 Regelin, E., Riede, H., Sandu, A., Taraborrelli, D., Tost, H. and Xie, Z.-Q., *The atmospheric*

1 *chemistry box model CAABA/MECCA-3.0*, *Geoscientific Model Development* **4**, 373–380,
2 2011.

3 Scalasca, <http://www.scalasca.org/>, Last Access: 23/11/2015.

4 Suarez, E., Eicker, N. and Gürich, W., *Dynamical Exascale Entry Platform: the DEEP*
5 *Project*, *inSiDE* **9**, 2, 2011.

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

1 Table 1: Experimental setup for JUDGE scalability test out-of-the-box.

Number of columns	8192 columns with 90 levels
Number of grid points	737280 grid points
Number of chemical species	139 species in 318 reactions
Spectral resolution	T42L90MA

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

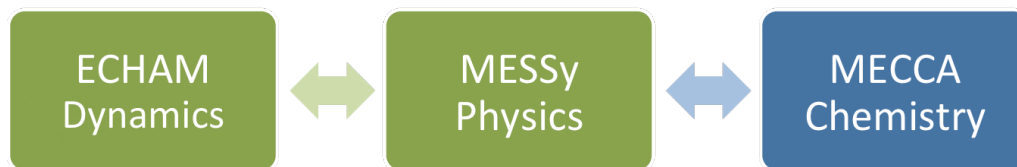
20

21

1 Table 2: System setup details for the analysis done on the JUDGE system.

Backend compiler version	Intel 13.1.3
MPI runtime version	Parastation/Intel MPI 5.0.27
Compilation flags	-O3 -fp-model source -r8 -align all
MPI processes per node	24

2

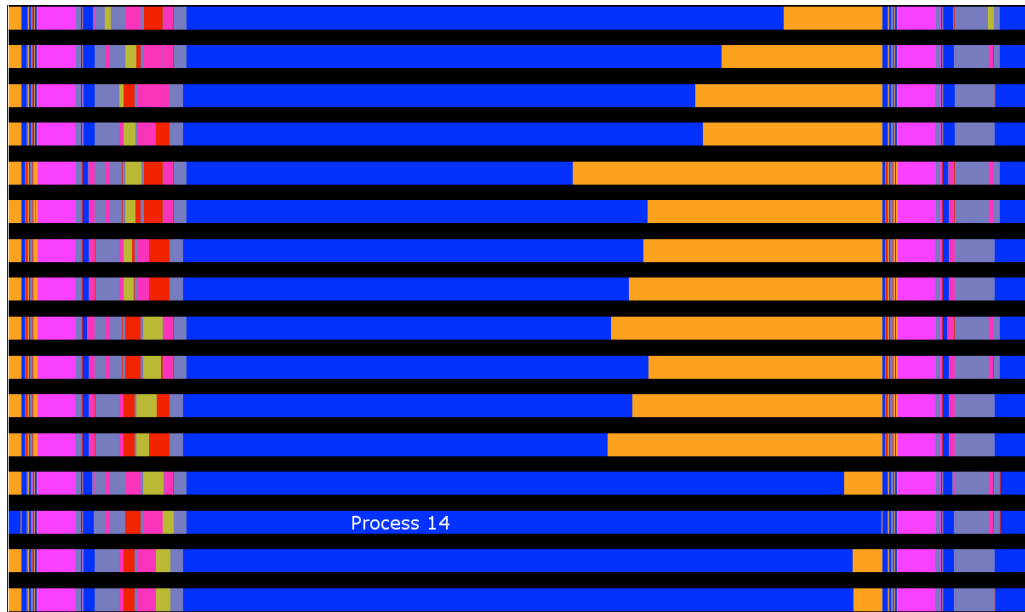


1

2 Figure 1 Phases of EMAC. Green phases run on the Cluster, blue phases run on the Booster.

3

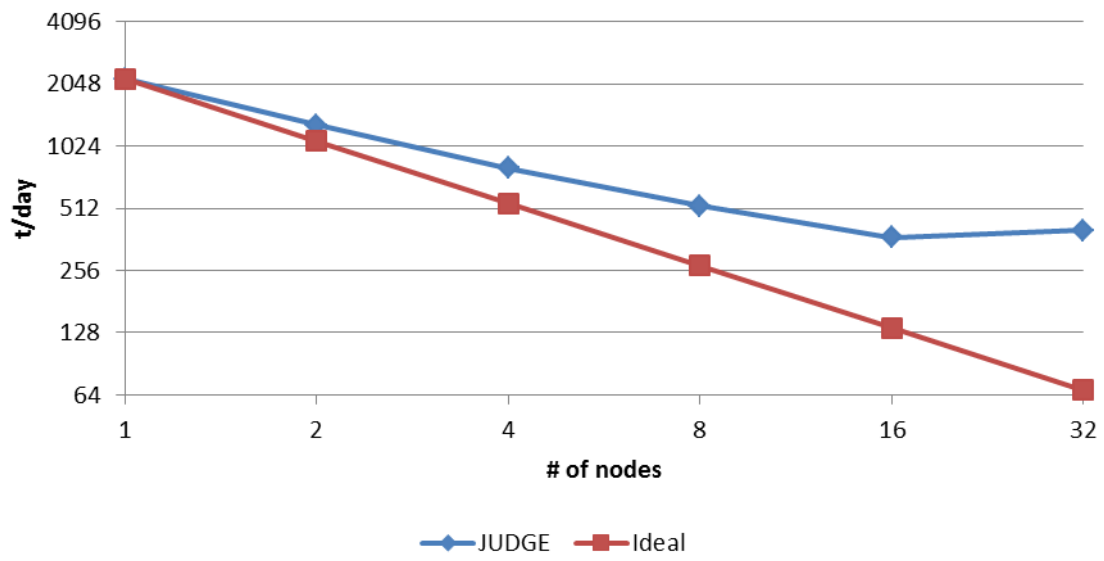
4



1
2
3
4
5

Figure 2 : Paraver trace of major processor usage of one time step. Time is along the horizontal and each bar corresponds to a separate CPU core. Blue colour depicts computation; orange corresponds to idle time due to load imbalance. The grid-space transpositions and Fourier and Legendre transformations are shown in magenta.

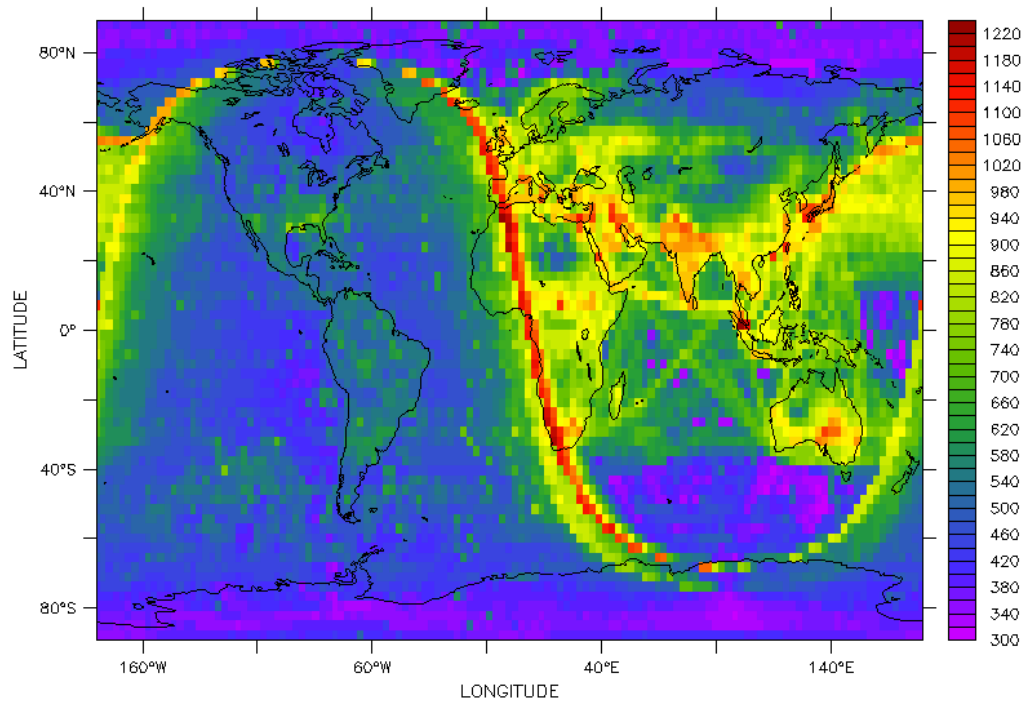
Performance of EMAC on JUDGE



1

2 Figure 3 : Wall time for one simulated day versus the number of nodes on JUDGE.

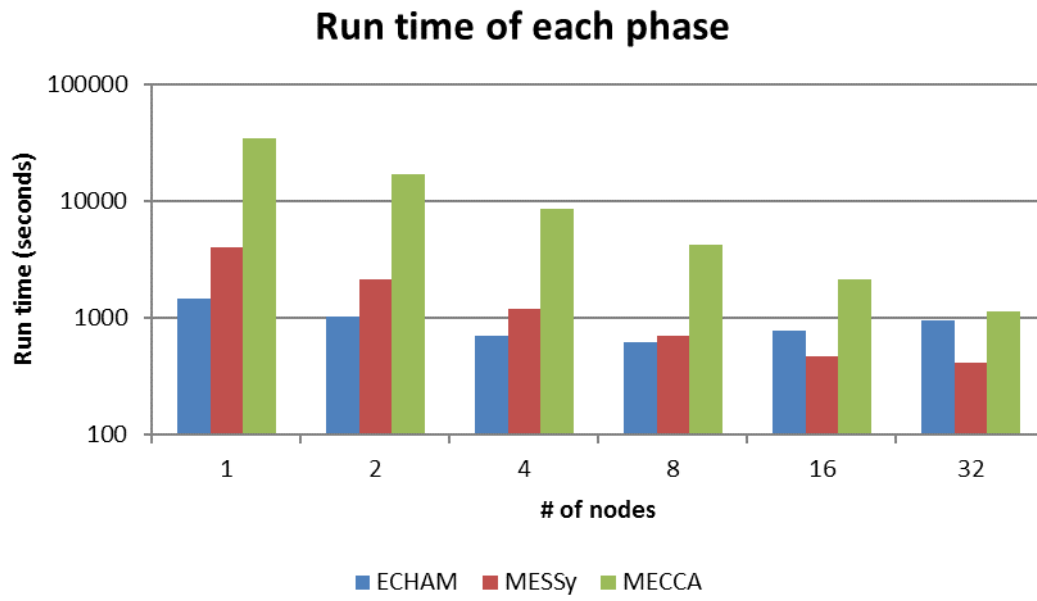
3



1

2 Figure 4 : Column-maximum MECCA kernel execution wall-time in microseconds for a single time-step. The
 3 adaptive time-step integrator shows a non-uniform run time caused by stratospheric photochemistry and natural
 4 and anthropogenic emissions.

5

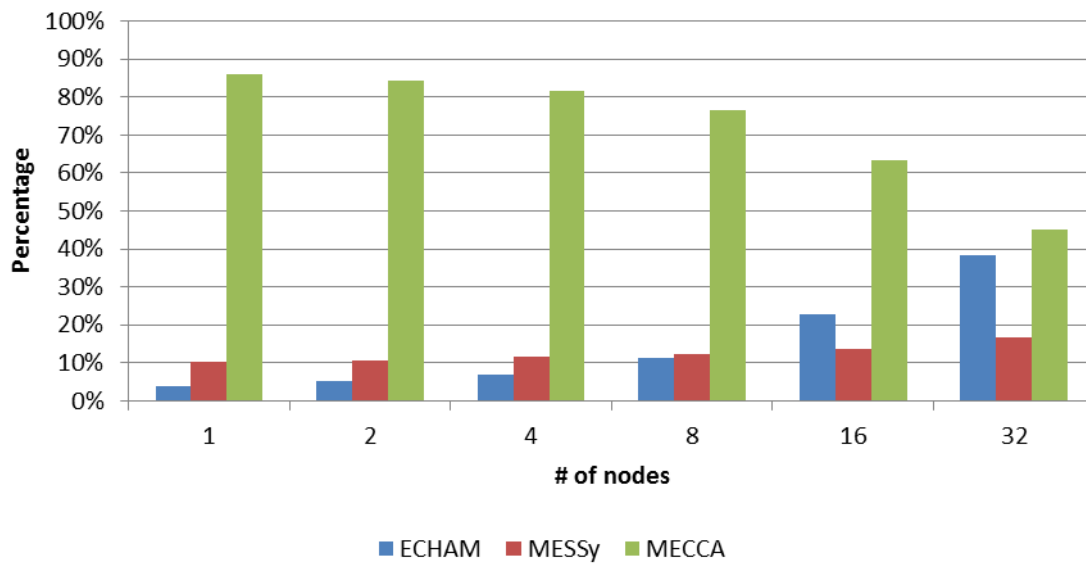


1

2 Figure 5 : Run time of each phase of EMAC, when running on MareNostrum 3.

3

Percentage of run time of each phase

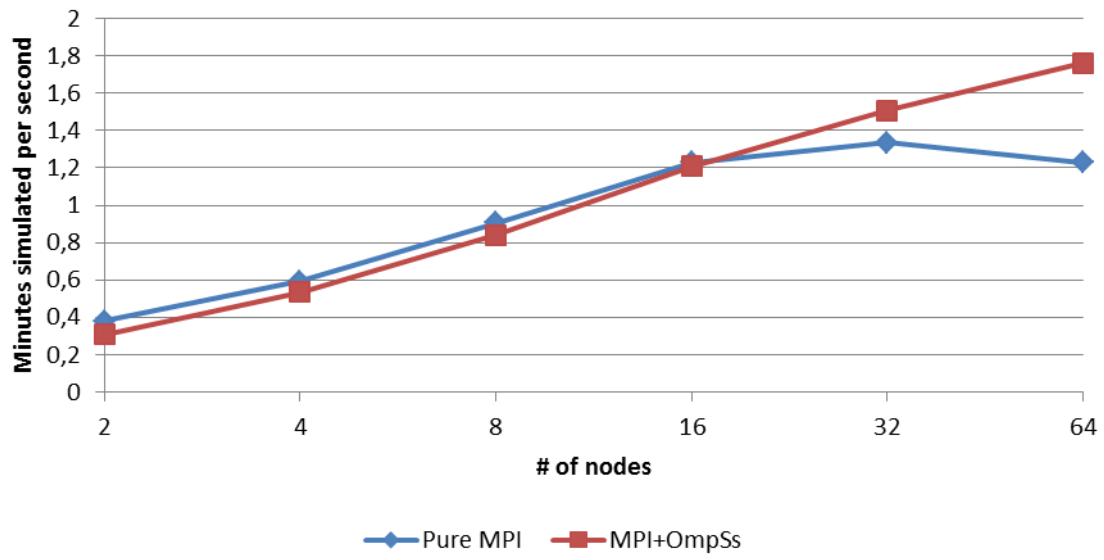


1

2 Figure 6 : Percentage of run time of each phase of EMAC, when running on MareNostrum 3.

3

Performance of OmpSs threading

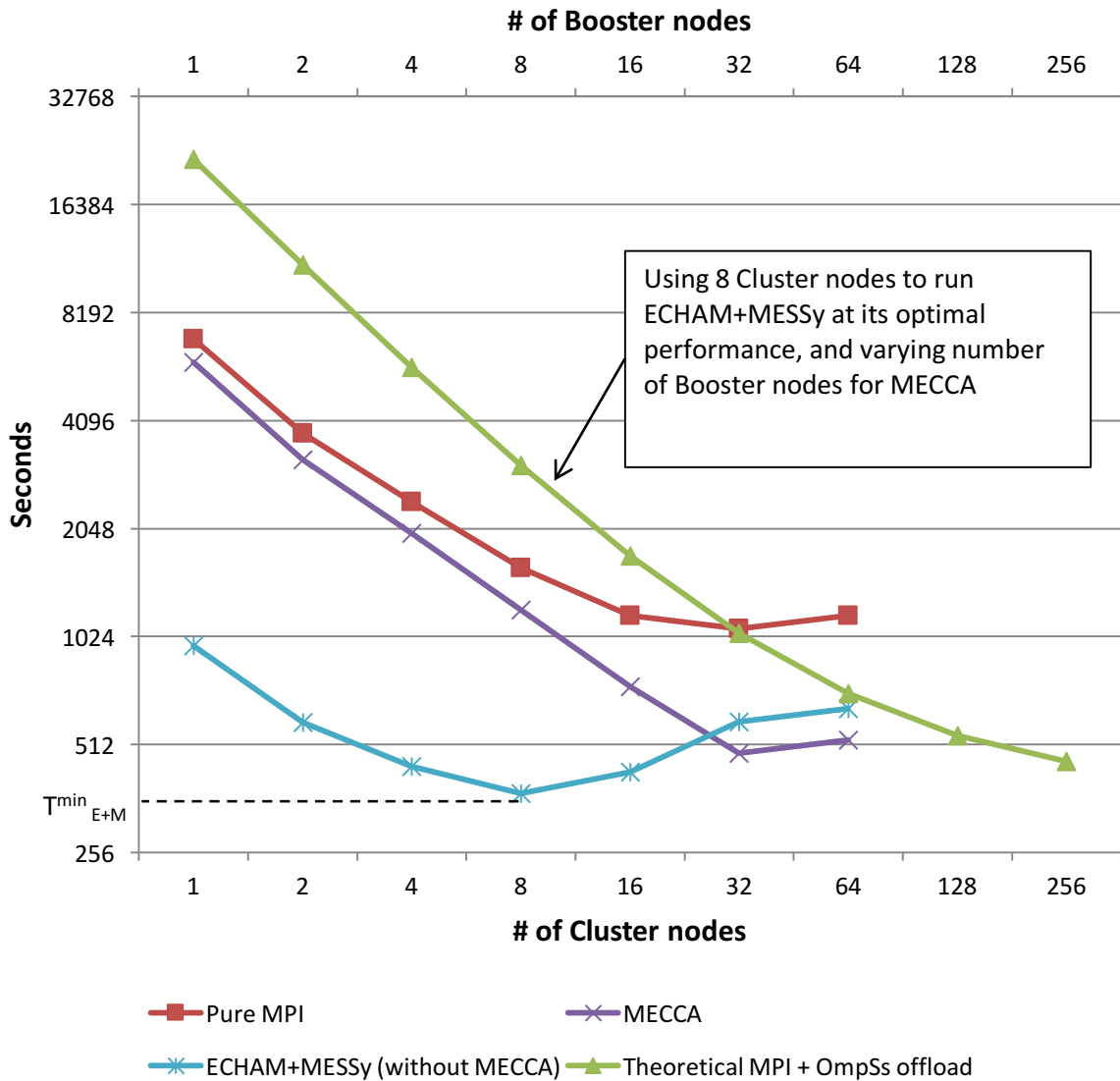


1

2 Figure 7 : Performance of OmpSs threading in the DEEP Cluster.

3

Theoretical performance with Xeon Phi offloading



1
 2 Figure 8: Time per simulated day in DEEP using a pure MPI approach, and a theoretical performance with
 3 offloading to Xeon Phi, based on the metrics collected in MareNostrum 3. The theoretical MPI + OmpSs offload
 4 data is based on a fixed configuration on the Cluster using 8 nodes, and scaling the number of Booster nodes.