Geoscientific
Model Development
Discussions

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

**M. Christou[1], T. Christoudias[1], J. Morillo[2], D. A. Mallon[3] 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.

Geoscientific
Model Development
Discussions

## 1   Introduction

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

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

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

The Cluster–Booster architecture is naturally suited to global atmospheric circulation–chemistry models, with global components running on the Cluster nodes exploiting the high-speed Xeon processors and local components running on the highly-parallel Xeon Phi co-processors. By balancing communication versus computation the DEEP concept provides a new degree of freedom allowing us to distribute the different components at their optimal parallelisation. The concept is depicted diagrammatically in Figure 1.

## 2  Overview of application structure

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

### 2.1  Phases

The implementation of EMAC comprises two main phases, the base model ECHAM integrating the dynamical state of the atmosphere, and the MESSy framework that interfaces to $n$ submodels calculating physical and chemical processes. Among these submodels stands out the MECCA submodel (Sander et al., 2007). This submodel computes the chemical kinetics of the homogeneous gas-phase chemistry of the atmosphere, and deserves special mention due to its intrinsic parallelism, high computational demands, and load imbalance rising from its rigid coupling to the base model's parallel decomposition.

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

It is, however, currently not possible to delegate the whole MESSy subsystem to full multi-threaded execution as some physical processes are naturally modelled in a column-based approach, and are strongly dependent on the system states at their vertically adjacent grid points. The implementation of submodels simulating these processes consequently relies on the column structure inherited from the base model. Furthermore, even a coarser column-oriented multi-threaded approach is hindered by global-variable interdependencies between submodels.

Describing homogeneous gas phase chemical kinetics, the MESSy submodel MECCA executes independently of its physical neighbours and is not limited by vertical adjacency relations. As more than half of the total run-time is spent in MECCA for a typical model scenario, it seems adequate to concentrate on the MECCA kernel with strong algorithmic locality and small

1  communication volume per task. As sketched in Figure **1** the current implementation of

2  MECCA, developed in the DEEP project, is delegated to the Booster using a task-based

3  approach while both ECHAM and the remaining MESSy submodels are executed on the Cluster

4  in the distributed-memory paradigm.

5  ## 2.2  Dominant factors

6  Implementing a spectral model of the dynamical state of the atmosphere, the ECHAM phase

7  comprises six transform and six transposition operations in each time step, as seen in Figure 3.

8  The data in memory for each time step (data size scales with the square of the model resolution)

9  is transposed in an all-to-all communication pattern, and this phase is dominated by network

10 bandwidth.

11 Figure 4 displays one time step traced with Extrae/Paraver (Extrae 2015, Paraver 2015) starting

12 with the end of the grid point calculations of the last time step—in which most processors are

13 already idle (orange) due to load imbalance and waiting for process 14 (blue) to finish running.

14 This is followed by the transpositions and Fourier and Legendre transformations (magenta),

15 which execute simultaneously as further analysis showed. After the transpositions a short

16 interval with all processors running (blue) can be identified with the time step integration in

17 spectral space, followed by the inverse transformations and transpositions and transport

18 calculations in ECHAM.

19 While the pattern described so far repeats towards the end of the displayed interval, the major

20 fraction of the time step is spent without communication, running (blue) or waiting (orange) in

21 calculations in MESSy in grid space. The MESSy phase comprises some 30 submodels that are

22 tightly coupled by exchanging the atmospheric observables using global variables. Model

23 performance depends largely on a virtual longitude run-time parameter exploiting cache line

24 adjacency of the grid point variables. Investigations during the first phase of the project

25 determined the load imbalance visible in Figure 4 to be caused by chemical processes computed

26 in the MECCA submodel.

27 The observed load imbalance is one of the main factors determining application scalability. It

28 is caused by an adaptive time-step integrator solving a system of differential equations. As the

29 stiffness of these equations representing homogeneous photochemical reactions varies by up to

30 two orders of magnitude due to changes in the intensity of sunlight, the adaptive integrator

31 demands varying amounts of run time accordingly (described in more detail in Section 2.3).

1 In the MECCA phase the algorithmically complex adaptive time-step differential equation

2 integrator operates on chemical concentrations of a total data size of the order of a few kilobytes

3 per grid point. Yet, as seen in Figure 5 that highlights the load imbalance caused by MECCA

4 and observed using Scalasca (Scalasca 2015), this phase consumes the major proportion of the

5 total execution time, it is compute-bound and an obvious candidate for offloading to

6 accelerators. It should be noted though, that in a regular architecture, accelerating this highly

7 parallel phase will not eliminate the load imbalance.

8 ## 2.3  Scalability considerations

9 To test the model scalability out-of-the-box, the EMAC application has been ported to the

10 JUDGE cluster at JSC, and a representative benchmark with a horizontal resolution of 128 grid

11 points in longitudinal and 64 grid points in latitudinal direction with 90 vertical levels and a

12 spin-up period of 8 simulated months has been compiled, frozen and packaged to be used for

13 measurements. Table **1** details the experimental setup for the results shown in this section.

14 EMAC was benchmarked with different numbers of processors on JUDGE in order to

15 determine the run time behaviour of the total application. As shown in Figure 6 the application

16 scales up to 384 processes (16 nodes x 24 MPI processes each), at higher numbers the

17 performance decreases. Parallel execution speed is determined by the balance of three factors:

18 computation, communication, and load imbalance. The benchmarking setup for the JUDGE

19 cluster can be seen in Table 2. While the computational resources increase with additional

20 processors and therefore increase the application performance, communication demands

21 diminish the positive effect of the additional processors. Additionally, increasing the granularity

22 of the total workload also increases the load imbalance.

23 While the number of processors used for the distributed-memory part of the code is limited by

24 the scalability of the non-local representation of global physical processes in ECHAM, the local

25 processes in MESSy running independently from their neighbours scale very well. The MESSy

26 subsystem has not been designed for multi-threaded execution, though, and contains non-local

27 code due to characteristics of the physical processes and algorithmic design decisions. Some

28 physical processes are naturally modelled in a column-based approach, because they are

29 strongly dependent on the system states at vertically adjacent grid points, e.g. sunlight intensity

30 at lower grid points depending on the absorption at higher grid points, and precipitation

Geoscientific
Model Development
Discussions

1    depending on the flux of moisture from vertically adjacent grid cells. Sub-models simulating

2    these processes consequently rely on the column structure implemented in the current model.

3    In the existing distributed-memory parallel decomposition, the three-dimensional model grid is

4    split horizontally using two run-time parameters, setting the number of processes in latitudinal

5    and longitudinal direction. As work is distributed independently for each direction, a

6    rectangular decomposition is obtained.

7    The physical load-imbalance, caused by photo-chemical processes in the lower stratosphere and

8    natural and anthropogenic emissions, appears in the run time spent for each grid point when

9    examining the benchmark calculations. In Figure 7 the maximal MECCA kernel execution

10   wall-time for one grid point in each column differs by up to a factor of four. The load imbalance

11   is caused by the adaptive time-step integrator solving the differential equations that describe

12   the chemical equations computed in the MECCA submodel. The strongly varying light intensity

13   at sunrise and sunset and night-time emissions lead to stiff differential equations that require

14   more intermediate time steps with derivative function evaluations and increase the

15   computational load by up to one order of magnitude.

16   At high levels of parallelisation, the load imbalance becomes a limiting factor, and the factors

17   determining scalability in absolute numbers in Figure 8 are both communication and

18   computation. For the ECHAM phase (blue), when scaling to beyond 8 nodes the

19   communication demands of the underlying spectral model involving several all-to-all

20   communication patterns start to dominate.

21   In Figure 9 the point at which communication and computation require equal times around 8

22   nodes is clearly apparent; 16 nodes is commonly used in production runs of the EMAC

23   atmospheric model as a scientific application to balance efficiency and total required wall time.

24   **3   Model Developments**

25   **3.1   Intranode taskification**

26   The EMAC model atmospheric chemistry code (MECCA) was taskified using OmpSs (Bueno,

27   J. et al., 2011, 2012, Duran, A. et al., 2011, Florentino et al., 2014) directives. OmpSs allows

28   the user to specify inputs and outputs for blocks of code or functions, giving enough information

29   to the runtime to construct a dependency graph. This dependency graph reflects at all moments

30   which tasks are ready to be executed concurrently, and therefore the programmer does not have

1  to explicitly manage the parallelisation. This idea of tasks and task dependencies has been
2  adopted in the OpenMP 4.0 standard (OPENMP 4.0, 2013). Since in MECCA each gridpoint is
3  completely independent of its neighbours, this part of the code is in principle embarrassingly
4  parallel, with no communication or inter-task dependencies involved.

5  The MECCA submodel was refactored through the creation of computational kernels for
6  intranode parallelisation with shared-memory tasks. The new version of EMAC, running
7  ECHAM with MPI processes and MECCA with shared-memory OmpSs tasks outperforms the
8  old EMAC using pure MPI, and continues to scale beyond the region where the original
9  implementation scaling performance plateaus. This can be seen in Figure 10, which shows the
10  performance using multi-threading on the DEEP Cluster.

11  **3.2   Internode taskification**

12  In DEEP, OmpSs has been extended to support offloading tasks to remote nodes (Beltran et al.,
13  2015). This mimics the behaviour of other accelerator APIs that move data from the host to the
14  device, compute in the device, and return the results to the host. However, OmpSs adds two
15  very important features: i) it allows offloading to remote nodes, not just locally available
16  coprocessors/accelerators, which is a key functionality to effectively use the Booster; and ii) it
17  allows using the Booster as a pool of coprocessors, so tasks can be offloaded to any Booster
18  node with enough free cores. The latter enables to eliminate the load-imbalance caused by
19  sunlight gradients in MECCA.

20  In a shared-memory taskification the data is already shared between threads, and no memory
21  copies are necessary. However, in DEEP, to leverage the Booster, this data has to be copied to
22  the Booster nodes. Keeping that in mind, the new task-based MECCA implementation was
23  optimised and the memory and network footprint of the distributed-memory offloading was
24  reduced by three orders of magnitude. To minimise the memory footprint for offloaded tasks,
25  the number of computational grid elements issued to MESSy is further split into individual
26  elements for each task, by rearranging the grid point arrays in each time step to implement data
27  locality at the grid-point level, resulting in a reduction of the total memory footprint from 2.7
28  MB down to 6.3 KB for each task. This was the result of refactoring both the data and code
29  structures in MECCA. At the benchmark resolution of T42L90MA a total number of 737 280
30  tasks are generated in each time.

1 As discussed in section 2.2, a detailed analysis of the EMAC run-time behaviour using Scalasca

2 (Wolf et al., 2008) and Extrae/Paraver has identified that the MECCA submodel consumes a

3 major proportion of the execution time, does not participate in communication, and is

4 independent of adjacency constraints. It is thus well suited to be delegated to the Booster

5 employing the large dynamical pool of accelerator resources provided by the DEEP concept for

6 load balancing of the heavily varying computation demands discussed in section 2.3.

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. Thus, the number of tasks to be sent to the Booster can be controlled and

11 optimised for each architecture, and host-specific configuration allows for optimum task size

12 based on bandwidth, reducing task communication overheads.

13 With this approach, the specifics of the DEEP system architecture, and in particular the

14 hardware present in MIC coprocessors is exploited by massively parallelising the chemistry

15 calculations at the gridpoint level and offloading to the Booster exposing a significant amount

16 of thread parallelism. At the same time the load imbalance observed in MECCA is

17 automatically alleviated through OmpSs' dynamic load balancing by selecting a sufficiently

18 fine task size and decoupling the model-domain location of the grid point from the task

19 execution on the physical CPU.

20 ### 3.3 Vectorisation

21 The computational core of MECCA is connected by an interface layer to the MESSy

22 framework, integrating different submodel code and data structures into the ECHAM base

23 model. It provides the gridpoint data as sub-arrays of the global simulation data – which have

24 been rearranged from their native longitude and latitude coordinates into a virtual longitude and

25 an outer index variable counting the virtual longitude blocks and assuming the role of a virtual

26 latitude. The virtual longitude exploits cache line adjacency on non-vector architectures and

27 serves as run-time vectorisation parameter for all MESSy submodels.

28 For the MECCA submodel an integrator kernel has been created that can be offloaded onto

29 worker threads running on the main processor or hardware accelerators. The chemical

30 mechanism is compiled by the Kinetic Pre-processor (KPP, Damian et al., 2002) implementing

31 a domain-specific language for chemical kinetics. The integrator kernel operates on the

1   variables of one grid-point describing the local state of the atmosphere and the integrator

2   parameters determining the solution of the chemical equations. As the kernel variables are

3   passed as one-dimensional sub-arrays of global, four-dimensional arrays, extending along the

4   virtual longitude, the vector variables are transposed to extend contiguously along the

5   dimension of chemical species.

6   In order to estimate the run-time effect of the changes in the code, the application was

7   benchmarked on the DEEP Cluster using the Xeon main processors without vectorisation as

8   baseline measurement of 548.65 seconds per simulated day. Compiling with the auto-vectoriser

9   enabled for the AVX instruction set extensions decreased the run time to 466.40 seconds,

10  resulting in a first speed-up of 1.18. Examination of the optimisation report identified several

11  unaligned array accesses, which were solved using compiler directives and introducing aligned

12  leading dimensions at 64-byte boundaries for multi-dimensional arrays as needed for the

13  instruction set of Intel Xeon Phi. These changes improved the total application performance to

14  349.50 seconds per simulated day for a second speed-up of 1.33 achieving a total speed-up of

15  1.57 (Figure 11).

16

## 4   Attainable Performance

18  At the time of writing this manuscript the DEEP Booster is in the bring-up phase, and not

19  available to users. In order to project the performance of the full DEEP System, Xeon-based

20  measurements on the DEEP Cluster were combined with Xeon Phi-based measurements on

21  MareNostrum 3. The DEEP Cluster reference data weighted by the relative factors for each

22  phase derived from the metrics measurements exhibit a performance maximum for the base

23  model (ECHAM) and MESSy (excluding MECCA) at 8 nodes, representing a good estimate

24  for the optimal parallelisation of that phase on the Cluster. This estimate of 375 s per simulated

25  day for the low-scaling Cluster phases was used to extrapolate the attainable performance,

26  merging this result at 8 nodes, with the Xeon Phi data retrieved from MareNostrum 3, where

27  benchmarks using one node had been run with varying numbers of processing elements within

28  one Xeon Phi processor. The pure MPI time in the DEEP Cluster, the time for each phase, and

29  the theoretical performance when offloading to the Booster are shown in Figure 12.

30  While the number of Booster nodes required to attain similar performance to the original

31  distributed-memory based implementation corresponds to regular accelerator architectures with

32  individual boosters directly attached to cluster nodes, the projected DEEP performance scales

beyond the optimal performance achieved so far. The EMAC atmospheric chemistry global climate model seems therefore well suited to exploit an architecture providing considerable more hardware acceleration than provided by regular systems. The projected attainable performance that outperforms the pure-MPI conventional cluster paradigm at higher core count (depicted here as the number of Booster nodes, while keeping the ECHAM/MESSy MPI part on 8 Cluster nodes for optimal performance) is also shown in Figure **12**.

## 5   Conclusions

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

The chemistry mechanism was taskified at the individual gridpoint level using OmpSs directives. The chemistry code was refactored to allow for memory adjacency of vector elements. Enabling the vectoriser achieves a total speed-up of 1.57 by aligning all arrays at 64-byte boundaries. The OmpSs taskification with remote offload allows for massive task parallelisation and the implementation of optional two-stage offload to control Cluster-Booster task memory size and optimum bandwidth utilisation.

The computational load imbalance arising from a photochemical imbalance is alleviated at moderate parallelisation by assigning grid points with differing run times to each process and distributing the load over all processes. Due to the physical distribution of sunlight this load balancing does not require an explicit algorithm at moderate parallelisation; instead, the implicit assignment of the model grid in rectangular blocks suffices for this purpose. At higher numbers of processors this implicit load-balancing decreases and the resulting load imbalance has to be solved by active balancing. The dynamic scheduling provided by the OmpSs run-time system balances the computational load without a possible, but expensive prediction for the current time step.

With these approaches, the specifics of the DEEP System architecture, and in particular the hardware present in MIC coprocessors can be exploited by massively parallelising the chemistry calculations at the gridpoint level and offloading to the Booster exposing a significant amount of thread parallelism. At the same time the load imbalance observed in MECCA will be automatically alleviated through dynamic load balancing by selecting a sufficiently fine task size and decoupling the model-domain location of the cell from the task execution on the physical CPU.

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

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

**Code Availability**

The Modular Earth Submodel System (MESSy) is continuously further developed and applied by a consortium of institutions. The usage of MESSy and access to the source code is licenced to all affiliates of institutions which are members of the MESSy Consortium. Institutions can become a member of the MESSy Consortium by signing the MESSy Memorandum of Understanding. More information can be found on the MESSy Consortium Website (http://www.messy-interface.org).

**Acknowledgements**

Geoscientific
Model Development
Discussions

# 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.

Geoscientific
Model Development
Discussions

1  Florentino, S., Mateo, S., Beltran, V., Bosque, J., L., Martorell, X. and Ayguadé, E.,

2  *Leveraging OmpSs to Exploit Hardware Accelerators,* International Symposium on Computer

3  Architecture and High Performance Computing, 112–119., 2014.


4  Jöckel, P., Kerkweg, A., Pozzer, A., Sander, R., Tost, H., Riede, H., Baumgaertner, A., Gromov

5  S. and Kern, B., *Development cycle 2 of the Modular Earth Submodel System (MESSy2),*

6  Geoscientific Model Development **3**, 717–752, 2010.


7  Mallon, A, D., Lippert, T., Beltran, V., Affinito, F., Jaure, S., Merx, H., Labarta, J.,

8  Staffelbach, G, Suarez, E. and Eicker, N., *Programming Model and Application Porting to the*

9  *Dynamical Exascale Entry Platform (DEEP),* Proceedings of the Exascale Applications and

10  Software Conference, Edinburgh, Scotland, UK., 2013.


11  Mallon, A., D., Eicker, N., Innocenti, M.E., Lapenta, G., Lippert, T. and Suarez, E., *On the*

12  *Scalability of the Cluster-Booster Concept: a Critical Assessment of the DEEP Architecture,*

13  FutureHPC '12, Proceedings of the Future HPC Systems: the Challenges of Power-

14  Constrained Performance, ACM New York, 2012.


15  OpenMP Application Program Interface Version 4.0 – July 2013, OpenMP Architecture

16  Review Board, http://www.openmp.org/mp-documents/OpenMP4.0.0.pdf, Last Access:

17  23/11/2015.

18  Paraver, Barcelona Supercomputing Center, https://www.bsc.es/computer-

19  sciences/performance-tools/paraver, Last Access: 23/11/2015.

20  Roeckner, E., Brokopf, R., Esch, M., Giorgetta, M., Hagemann, S., Kornblueh, L., Manzini,

21  E., Schlese, U. and Schulzweida, U., *Sensitivity of simulated climate to horizontal and*

22  *vertical resolution in the ECHAM5 atmosphere model,* J. Climate **19**, 3771–3791, 2006.

23  Roeckner, E., Bäuml, G., Bonaventura, L., Brokopf, R., Esch, M., Giorgetta, M., Hagemann,

24  S., Kirchner, I., Kornblueh, L., Manzini, E., Rhodin, A., Schlese, U., Schulzweida, U. and

25  Tompkins, A., *The atmospheric general circulation model ECHAM 5. PART I: Model*

26  *description,* Report/MPI für Meteorologie **349**, 2003.

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.

| Scaling | Number of columns | Number of grid points | Number of chemical species | Spectral resolution |
|---|---|---|---|---|
| Strong scaling | 8192 columns with 90 levels | 737280 grid points | 139 species in 318 reactions | T42L90MA with 42 coefficients |

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

2

1



2  Figure 1 : Distribution of the Earth System Model components on the Cluster-Booster architecture.

3

4

1

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

3

4

Geoscientific
Model Development
Discussions

Open Access

EGU



1

2   Figure 3 : The ECHAM main application loop. MESSy replaces and enhances the grid point calculations marked

3   in yellow.

4

5

1

2  Figure 4 : Paraver trace of major processor usage of one time step. Time is along the horizontal and each bar

3  corresponds to a separate CPU core. Blue colour depicts computation; orange corresponds to idle time due to load

4  imbalance. The grid-space transpositions and Fourier and Legendre transformations are shown in magenta.

5

Geoscientific
Model Development
Discussions

1



2    Figure 5 : MECCA execution time analysed with Scalasca.

3

Figure 6 : Wall time for one simulated day versus the number of nodes on JUDGE.

Geoscientific
Model Development
Discussions
Open Access

EGU



1

Figure 7 : Maximal MECCA kernel execution wall-time in microseconds. The adaptive time-step integrator shows

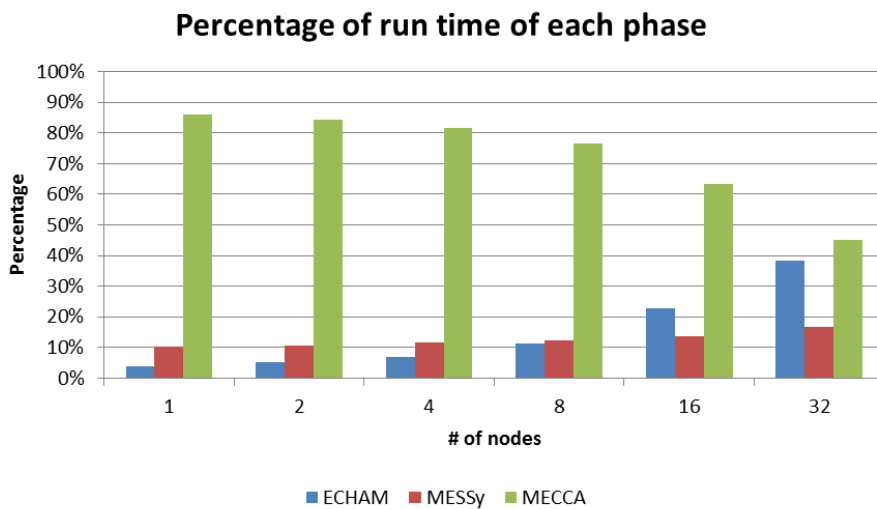a non-uniform run time caused by stratospheric photochemistry and natural and anthropogenic emissions.

1

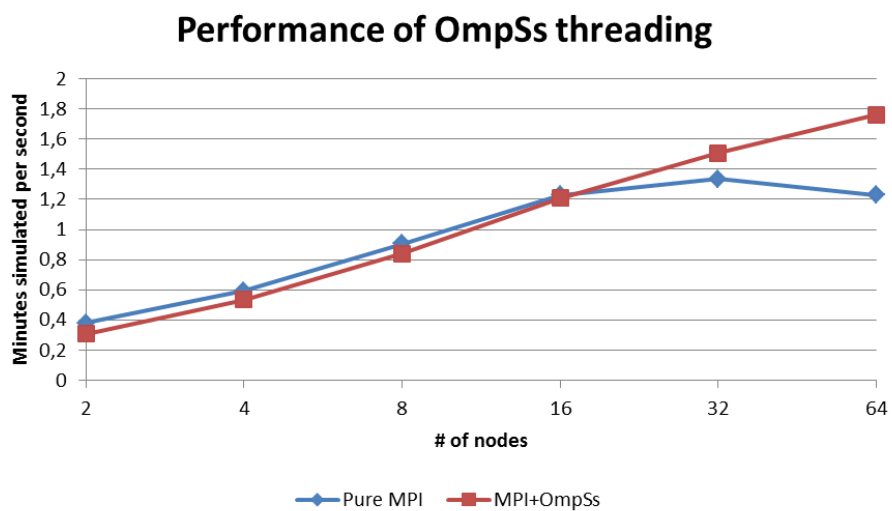2    Figure 8 : Impact on run time of each phase of EMAC, when running on MareNostrum 3.

3

1

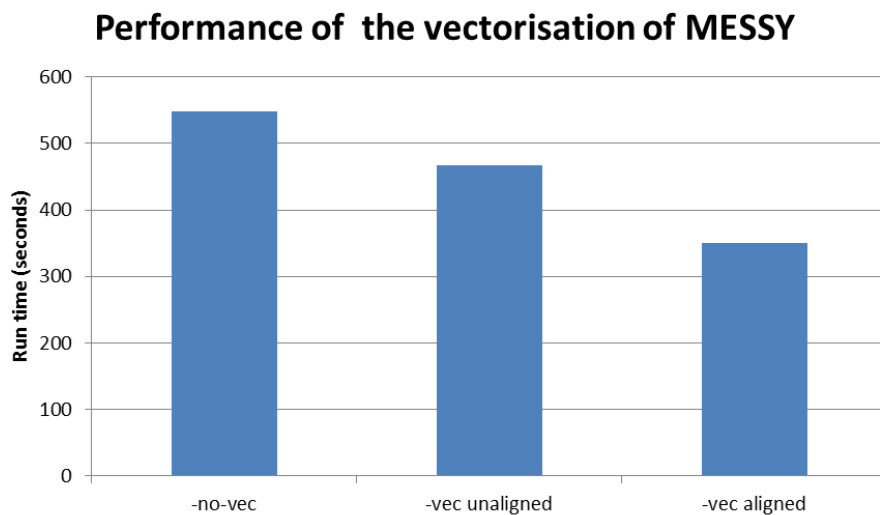2 Figure 9 : Percentage of run time of each phase of EMAC, when running on MareNostrum 3.

3

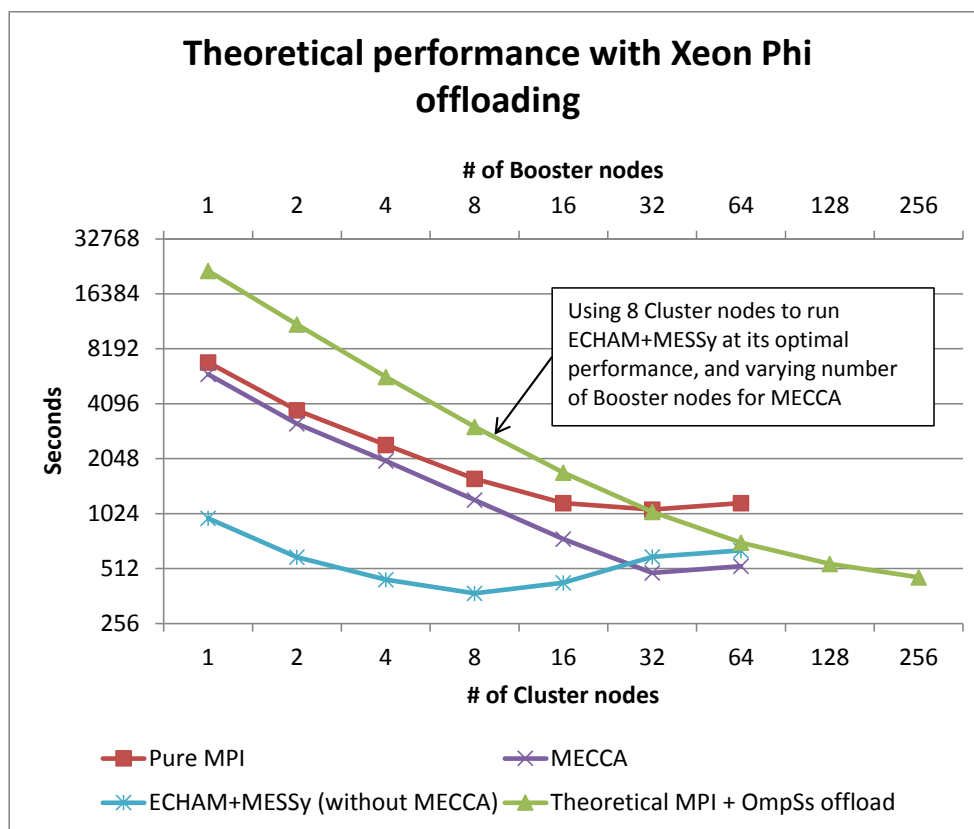Figure 10 : Performance of OmpSs threading in the DEEP Cluster.

Figure 11 : Performance of the vectorisation of MESSY in a Xeon E5-2680.

Figure 12: Time per simulated day in DEEP using a pure MPI approach, and a theoretical performance with offloading to Xeon Phi, based on the metrics collected in MareNostrum 3. The theoretical MPI + OmpSs offload data is based on a fixed configuration on the Cluster using 8 nodes, and scaling the number of Booster nodes.