

Effect of MPI tasks location on cluster throughput using NAS

Manuel Rodríguez-Pascual · José A. Morínigo · Rafael Mayo-García

Received: date / Accepted: date

Abstract In this work the Numerical Aerodynamic Simulation (NAS) benchmarks have been executed in a systematic way on two clusters of rather different architectures and CPUs, to identify dependencies between MPI tasks mapping and the speedup or resource occupation. To this respect, series of experiments with the NAS kernels have been designed to take into account the context complexity when running scientific applications on HPC environments (CPU, I/O or memory-bound, execution time, degree of parallelism, dedicated computational resources, strong- and weak-scaling behaviour, to cite some). This context includes scheduling decisions, which have a great influence on the performance of the applications, making difficult to achieve an optimal exploitation with cost-effective strategies of the HPC resources. An analysis on how task grouping strategies under various cluster setups drive the execution time of jobs and the infrastructure throughput is provided. As a result, criteria for cluster setup arise linked to maximize performance of individual jobs, total cluster throughput or achieving better scheduling. To this respect, a criterion for execution decisions is suggested. This work is expected to be of interest on the design of scheduling policies and useful to HPC administrators.

Keywords MPI application performance · Benchmarking · Cluster throughput · NAS Parallel Benchmarks

1 Introduction

The evolution in processors during the last decade has been oriented towards an increasing degree of parallelism and this fact has deeply impacted all levels

Manuel Rodríguez-Pascual · José A. Morínigo · Rafael Mayo-García
Departament of Technology, Centro de Investigaciones Energéticas, Medioambientales y Tecnológicas, CIEMAT, Madrid, Spain
Corresponding author: E-mail: josea.morinigo@ciemat.es

of computing hardware and software. The design of clusters and supercomputers is also following this path. For example, the number of cores according to the TOP500 list [14] has grown exponentially since 1993. Current trends include the use of many-core processors, driving the number of computing units even further. An obvious way of getting the most of this is the usage of highly parallel applications. MPI and OpenMP continue being instrumental to create highly scalable applications suitable for this environment. Applications can be classified into CPU, I/O or memory-bound depending on which factor limits its execution speed. Other common issues regarding requirements are execution time, degree of parallelism and required computational resources, to cite some. This leads to many questions and specific scheduling decisions. Probably the first is what should we do with partially-filled multi-core processors? What is the most efficient decision when a given job is only using some of the CPUs of a node, or just some of the cores of a CPU? On one side, executing some other job at the same time would lead to a most intensive usage of the resources; on the other, sharing the resources (namely memory at different levels and I/O) may force both jobs to compete for them and slow down, having in fact a negative impact on the total execution time. Thus, the implementation of right task grouping strategies managed by the cluster scheduler seems to be a fundamental aspect to speedup the executions by setting the optimal location of the job MPI tasks. The problem is complex and has extra dimensions, and the scope of this work aims to shed some light on performance issues.

2 Related work

Impact of MPI task locality has been investigated in [7] with three kernels of NAS Parallel Benchmarks (in short, NPB) and application codes running in a cluster of 32 CPUs. They show that an execution time saving of up to 25% is possible. Their number of used CPUs is small and the authors plan to extend the experiments to large-scale machines since it seems necessary to be conclusive about what happens in situations of many processors. In [4] it is summarized the results of mapping MPI tasks onto sockets, taking into account the machine topology. Results show that it is beneficial to map tasks onto as many sockets per node as possible (the bigger savings in execution time, up to 30%, are obtained precisely for those cases). Similar experiments are done in [13], reporting an improvement of about 15%. In particular, research dealing with multicore architectures has been focused in the last years. To this regard, [11] presents the gain in computational efficiency of a MPI-based production application that exhibits a performance peak improvement of about 9% (with averaged performance improvement of 6%), attributed to a better use of cache-sharing at the same node and to the high intra- to internode communication ratio of the cluster. Although it is a modest speedup, it is noticed that it is obtained with minor source code modifications. The work in [3] points to the same direction by evaluating the impact of multicore architectures in a set of benchmarks; but on the contrary, they conduct a non-straightforward

adaptation of the original application. Their characterization of the inter- to intranode communications ratio throws a figure of 4 to 5 in the worst case. This kind of node mappings is an area where little to moderate efforts are required for significant gains in application performance. The impact of internode and intranode latency is analyzed in [12] using a parallel scientific application with MPI tasks mapped onto the CPUs of an infiniband-based cluster of 14 nodes. With the objective of improving the computational efficiency, [17] analyzes how many cores per node should be used for applications execution. Here, both NPB and a large-scale scientific application code are executed in three single-socket-per-node clusters. They identify that task mapping is an important factor on performance degradation, being the memory bandwidth per core the primary source of performance drop when increasing the number of cores per node that participate in the computation. Something similar concludes [10], showing a high sensitivity of the attained NAS kernels performance to the multi-core machines. In [16] it is detected that NPB exhibit high sensitivity to the cluster architecture. Also, MPI tasks mapping reveals that distributing them over the nodes is better from a computational standpoint in most cases. According to their experiments, an up to 120% speedup is attained for most of the NAS kernels. They explain this behaviour because by distributing the tasks they do not have to compete for node local resources, a scenario that seems to occur when running tasks are sharing a slot or are located in slot of the same node. In [15] a semi-empirical predictive model is formulated and tested with a large-scale scientific application code, which provides good results for weak scalability cases and show that it can lead to a 5% increase of the execution time. The study conducted in [5] on mapping MPI tasks to cores using micro-benchmarks and NPB, shows that it may affect significantly the performance of intranode communication, which is closely related to the inter- to intranode communication ratio.

These previous investigations point out to the large sensitivity of the execution time to the task mapping. The impact of grouping or not MPI tasks outside of the box (out of the same node), over sockets or cores within the node is high as it is seen that execution time varies significantly. Also the lack of understanding of how to proceed in a systematic manner with an application in a specific cluster remains and further clarifications are needed to improve the cluster efficiency and usage. The present investigation summarizes the results of mapping MPI tasks onto cores in two different infiniband-based clusters at CIEMAT.

Then, processor mapping combinations have been tested to explore the impact on cluster throughput and hence, to build usage criteria aiming at feeding better scheduling strategies to support the scientific groups.

Hence, the present work explores the behaviour of different scientific kernels on modern infrastructures and the impact on clusters throughput. An analysis on how task location, network traffic and resource sharing affect their execution time has been carried out to infer a generalization of their behavior. This information can then be useful to improve scheduling algorithms and cluster setups. In the text, a *job* is composed by one or more *tasks*. The assignment

of those tasks to computational resources (*mapping*) determining when and where to run each job constitutes the *scheduling* process.

3 Characterization of HPC facilities

3.1 Benchmarking

The chosen applications for systematic benchmarking can be divided in two groups. The first one is system benchmarks, to measure the raw performance of the components of our clusters. The second one is application benchmarks (NPB), to test the behavior of the infrastructure when running real applications.

STREAM benchmark [8] measures sustainable memory bandwidth and tests the communication bandwidth between the socket and its RAM. In multicore sockets an OpenMP flag is set for building one thread/core during compilation.

OSU micro-benchmark [9] measures the latency and bandwidth of MPI libraries and interconnections. It implements a set of routines with various communication patterns. It measures intra- and internode communication bandwidths.

Bonnie++ [2] is a small yet powerful benchmark to measure disk performance. It provides a number of tests of hard drive and file system performance.

Intel Memory Latency Checker 3.1 [6] allows accessing memory chunks located in the elements of the memory hierarchy, measuring the latency time.

All together, these benchmarks allow characterizing the cluster raw performance. Hence, a better understanding of the results gathered with a set of scientific kernels is intended. The NAS Parallel Benchmarks [1] is developed at the Numerical Aerodynamic Simulation (NAS) program at NASA. It has evolved as a group of kernel benchmarks, set for a variety of problem sizes (classes A, B, C,... which differ mainly in the sizes of the arrays) of increasing computing cost. All together are representative of algorithmic building blocks found in the aforementioned scientific kernels. Among them are examples of memory-bound and CPU-bound, or weak-scaling and strong-scaling kernels. The present investigation uses NPB v2.0, which includes seven portable kernels (Fortran90, MPI parallelised) whose acronyms corresponds to: BT- Block Tridiagonal solver; CG- Conjugate Gradient; EP- Embarrassingly Parallel; FT- Discrete fast Fourier Transform; IS- Integer Sort; LU- Lower-Upper Gauss-Seidel solver; MG- MultiGrid on a sequence of grids.

3.2 Infrastructure characterization

Two computing facilities of different generations based at CIEMAT have been employed (see Figure 1, which displays the hardware of their nodes). The first one, EULER, is a production HPC shared among the scientific groups. It

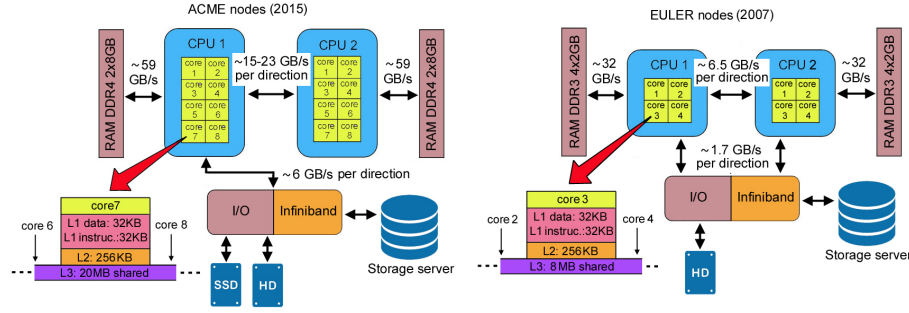


Fig. 1 Major architecture and communications for clusters ACME and EULER (bandwidths are approximate best cases sustained values).

consists in 480 Xeon®5450 quadcore@3.0 GHz, 2GB RAM/core, mounted on Dell PowerEdge M610 blades. When EULER was installed in Autumn 2008, its performance according to LINPACK (23 TFlops peak) would have ranked it around position 300 of TOP500. Because it is under full usage, the experiments have been done while sharing the resource and the restriction of accessing to a reserved set of 8 nodes within a limited timeframe (this matches with how the analyzed applications behave in real environments).

The second one, ACME, is a smaller, state-of-the-art HPC for research purposes and fully devoted to this project, which counts with 10 nodes (8 of them are computing nodes): 2 Bull R424-E4 chassis with 4 nodes each, plus another two devoted to host accelerators and fast network storage. Each node consists in a Supermicro X10DRT-P motherboard with two 8-core Xeon®E5-2640 v3@2.60. Each node counts with 32 GB DDR4 RAM memory @2133 MHz, in the form of 4 x 8 GB modules. Two of these modules are connected to each processor (NUMA). Each core accesses to half the memory with rather smaller access time than to the other half, as show in Fig. 2. Network is provided by Infiniband FDR. The MPI library mvapich2-2.2b is installed in both. It includes some performance metrics obtained with the benchmarks and from the official documentation.

Fig. 3 shows the dependence of internodes and intranodes communication bandwidth with the message size for the OSU micro-benchmark running on ACME and EULER clusters. It is worth noticing that the results mostly match their counterparts provided by the hardware vendors. ACME has higher intranode communication bandwidth (3 to 4 times higher than EULERs); 2.5 times larger shared L3 cache; and 3 times higher infiniband bandwidth. Both clusters have a ratio of intra- to internode bandwidth within 3 to 3.5.

3.3 Influence of node sharing on memory access time

In multi-core CPUs and in multi-CPU nodes, there is a decision concerning whether it is better or not to share the resources among pending jobs, or should a sole application be executed at the same time. It can be inferred

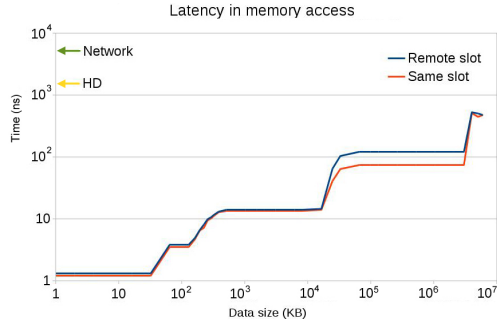


Fig. 2 Latency of a core accessing increasingly larger data blocks in ACME, corresponding to cache, RAM and disk. ‘Same slot’ and ‘Remote slot’ refer to the location of the processor to which the accessed memory is connected in the NUMA system.

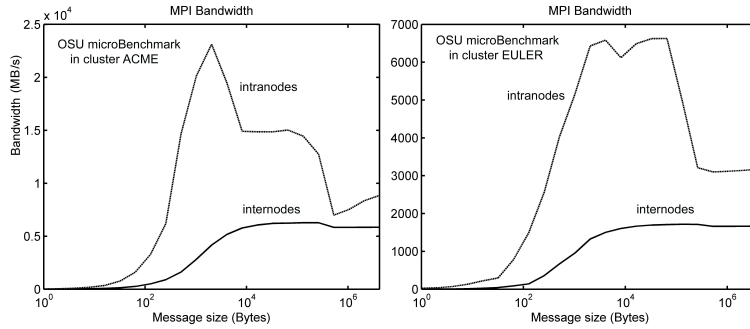


Fig. 3 Intranode and internode communication bandwidth in ACME and EULER clusters, computed with the OSU micro-benchmark.

that sharing a node between two or more jobs may increase RAM misses, as the available memory is shared between the running jobs, so they have less space for data and executable. This same issue happens when sharing a multi-core CPU, leading to an increased number of misses in L3 cache. Hence, it is necessary to measure the access time to the different memory layers in order to quantify its influence.

Figure 2 shows latency in ACME when a given core is accessing to memory. As expected, L1 cache is the fastest one but only stores up to 32 KB of data; then comes L2 with 128 KB, L3 with 20 MB, and after that the RAM memory. In this case, as pointed out before, there is a significant difference (60%) between accessing the modules connected to the same processor and those connected to the other one in the other slot of the same motherboard. The last step corresponds to the sizes between 32 and 64 GB, where both RAM modules are accessed to store/read. The last step of the memory hierarchy is represented by the persistent storage. It counts with a HD (1TB) for scratch and temporary files and network storage for the users home directory and the non-OS applications (scientific codes). Measured latency is $1.5\mu\text{s}$ for the HD

Table 1 Total execution time of NAS kernels in ACME with different cluster setups (time ratio is referred to the Dedicated Network setup)

Setup	Execution time (s)	Time ratio (%)
Dedicated Nodes	36533	32.6
Dedicated Cores	19442	17.5
Dedicated Sockets	28989	25.8

and $7.35\mu\text{s}$ for the network storage, roughly an order of magnitude larger than RAM latencies.

Summing up, results show that a miss in any cache level implies accessing an upper layer of the memory hierarchy with a penalization of about an order of magnitude in latency. Resource sharing will then have an impact on the job execution due to the influence on the access time of the different cache levels.

4 Results

In this section, results of executing a set of seven algorithms of the NAS Parallel Benchmarks are presented. These span a rather different behaviour in terms of CPU and memory requirements, execution time (ranging from a second to several hours) and have been executed in different ways to obtain as much information as possible regarding the machine performance, with the objective of identifying tasks execution patterns to extract general conclusions. The study comprises jobs up to 128 MPI tasks. This maximum number of tasks is justified by the limited strong-scaling properties of NAS kernels, as shown in Fig. 4, where it is visible that speedup stagnation occurs at about 32 MPI-tasks for most of them (stagnation for BT and EP begins later).

4.1 Cluster performance

The experiments with NPB have been repeated under four Slurm setups:

- Dedicated Cores: one-to-one assignment of cores to MPI tasks of the parallel job (a core executes only one MPI task of that job; set in ACME and EULER).
- Dedicated Sockets: a socket may only execute MPI tasks of the same job. Once the socket is occupied by at least one MPI task of a job, no other part of another job may be executed on it in the meanwhile (set in ACME).
- Dedicated Nodes: an entire node is assigned in exclusivity to execute MPI tasks of the same job (set in both ACME and EULER).
- Dedicated Network (reference case): the whole cluster executes only one parallel job at the same time, thus avoiding any overhead due to network congestion (set in ACME). This is a scenario devoted to obtain reference execution times.

Table 1 compares the total execution time of a set of NAS kernels. To mimic real-life workloads, all jobs corresponding to different kernel classes (sizes of

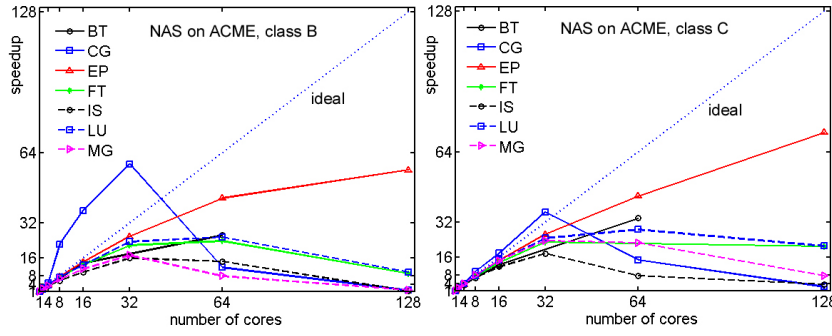


Fig. 4 Speedup of NAS kernels in cluster ACME. Intermediate size classes B (left) and C (right) are plotted.

Table 2 Submitted jobs of all NAS kernels partitioned per degree of parallelism.

Degree of parallelism	1	2	4	8	16	32	64	128
Total number of jobs	210	360	630	720	840	540	400	170
Percentage of jobs (%)	5.4	9.3	16.3	18.6	21.7	14	10.3	4.4

computed problems) and degrees of parallelism (about 4,000 jobs sent for each cluster setup) were submitted at the same time, letting Slurm scheduler to decide where and when to execute them using 8 nodes (16 cores/node) in ACME and 16 nodes (8 cores/node) in EULER. There was no indication of any job maximum execution time, thus no preemption techniques were employed. Table 2 shows the jobs according to their degree of parallelism; note that 128 is the number of cores in the cluster ACME. This way, the impact of MPI tasks location inside the clusters has been analysed under the Slurm setups. It is noted that the Dedicated Sockets setup is the 2nd more efficient after the Dedicated Cores setup (about 30% of the jobs counts for 8 or even less MPI tasks), which is able to allocate more than one job at the same time.

4.2 NAS Benchmarking

The number of nodes (nN) times the number of MPI tasks per node (nT), in short $nN \times nT$ (see Fig. 5), that define the configuration of each experiment (say, a definite kernel of a given class, executed under a cluster setup) has been repeated 10 times, with their average referred in what follows as a computed case. The experiments conducted under Dedicated Network setup include the kernels of class D to enlarge the population of computed cases. For the other cluster setups, only experiments with the A, B and C classes of the kernels have been conducted to guarantee that the running MPI processes fit into the RAM memory as well as to bound the workload computing time. All exhibit standard deviation $< 1\%$.

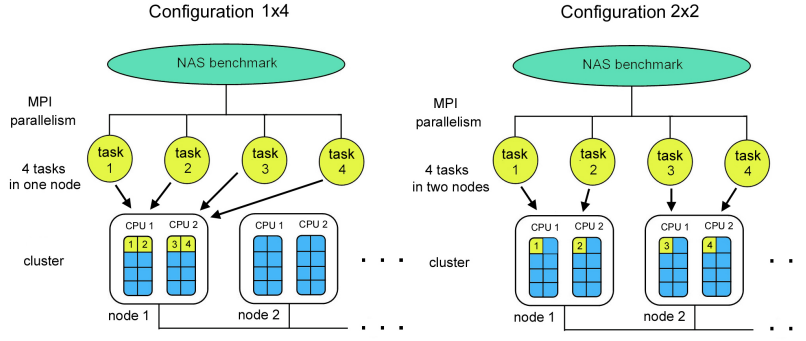


Fig. 5 Example of 1x4 and 2x2 MPI tasks mapping in cluster ACME.

4.3 Dedicated Nodes cluster setup

Bearing in mind that the Dedicated Cores setup is a realistic scenario of production clusters, a general trend can be stated out of Fig. 6, which depicts the execution time of the seven NAS kernels in ACME. Each sector of the circles corresponds to a given NAS kernel, identified by its acronym. Sectors are partitioned into annular regions: each annulus corresponds to a $nN \times nT$ configuration (indicated along the horizontal line on the right of the circle). That is, a circle annulus compares the impact of a prescribed $nN \times nT$ configuration on the execution time for all the NAS kernels executed under a degree of parallelization and the same NAS class. The outermost annulus corresponds to the most distributed scenario in terms of number of participating nodes (on the contrary, the circle centre corresponds to the most concentrated one). The execution time is shown dimensionless (referred to the execution time obtained in the cluster configuration of the fewest participating number of nodes, which corresponds to the circle centre). The examination of the relative execution times in Fig. 6 shows that many computed cases take more execution time as far as more nodes are involved.

Fig 7 compares ACME and Euler behaviour using the relative execution time (RET) as a performance metric. The plotted bars for each kernel gather the corresponding data of the NAS classes and $nN \times nT$ configurations (where the mean and maximum-minimum values are indicated), in order to see what is achieved in average. Both ACME setups -Dedicated Nodes and Dedicated Cores- show that execution time exhibits greater variation than EULER's (in particular, it is visible for the Dedicated Nodes setup in ACME, that the relative execution time may be quite large in some experiments: CG: 30.6; IS: 24.0). In addition, EULER exhibits more sensitiveness to the cluster setup. Fig 7 shows that all the mean values in ACME have $RET > 1$, thus it can be stated that grouping MPI tasks within as few nodes as possible is good to achieve shorter execution times. This statement changes in EULER, as it is visible in Fig 7b, which shows that most mean values have $RET < 1$ (and most of the kernels experiments with $RET < 1$: BT, FT, IS and MG). This

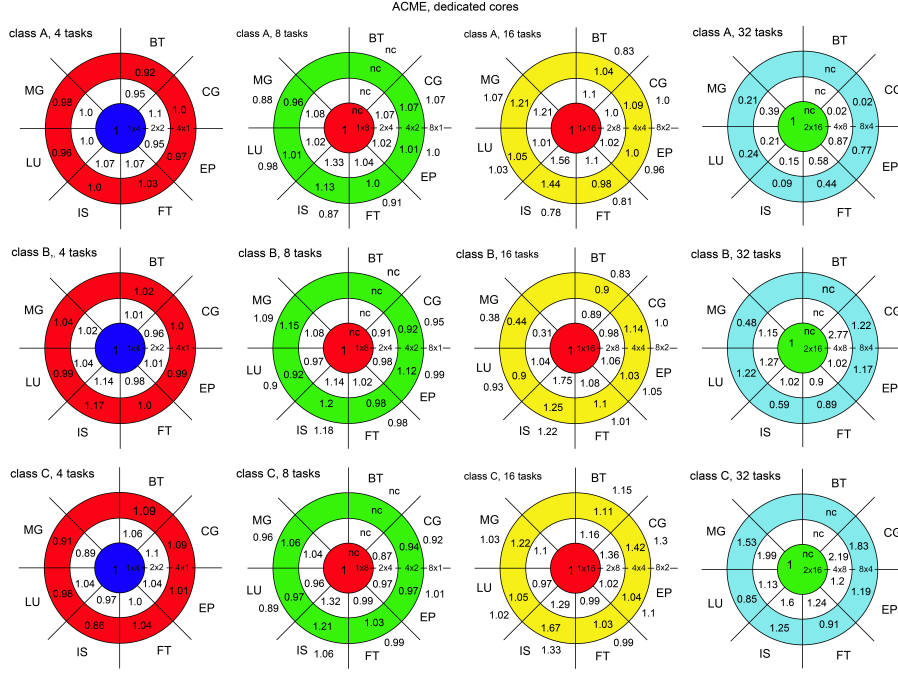


Fig. 6 Relative execution time for the Dedicated Cores setup in ACME. Nondimensional computing time is referred to the case of all processes running in one node.

result departs from ACME's. That is, Euler has many experiments in which the distribution of the MPI tasks over the nodes implies an execution speedup. The fewer host memory per core and slower communication bandwidth with the RAM stands as one cause to explain this trend. So, it seems that a general rule can be inferred out of the plots of both clusters after examining the behaviour of the kernels as a whole. However, several considerations must be pointed out regarding the computed behaviour. On one hand, a case by case examination reveals that there are exceptions to this rule as for example occurs in ACME regarding the LU kernel (class B - 8 tasks), IS kernel (class C - 4 tasks) and others. So non consistent tendencies of the kernels are noticed. And on the other hand, non monotone behaviour occurs in several computed cases. Some of these may be explained by the kernel requirements (CPU- or memory intensive,...), but also it is suggested that the execution is affected by the MPI tasks of rather different behaving kernels located in neighboring cores by the scheduler, which compete for resources (RAM and traffic). Besides, it is interesting to notice that such a pattern (and criterion) related to grouping MPI tasks in fewer nodes, seems to be more effective in saving execution time as the degree of parallelism and size (class) increase. On the contrary, computed cases of low number of tasks (see column of 4 MPI tasks in Fig. 6) show a small variation of the execution time (within 5-10%) with the number of nodes. As said, this general trend observed in ACME, it is not so definite in EULER. The

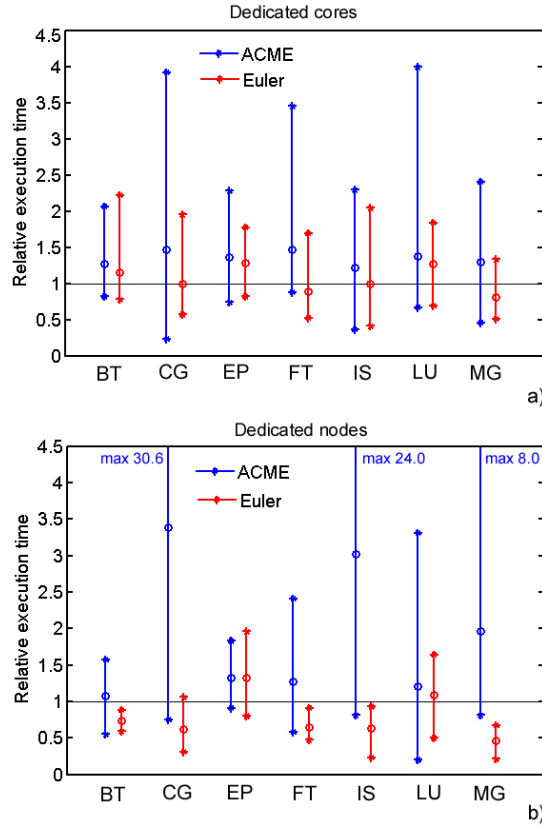


Fig. 7 Execution time of the NAS kernels for the setups Dedicated cores and Dedicated nodes in both clusters ACME and EULER. Times are nondimensionalized by the most concentrated case. Bars have been plotted using the entire set of data (NAS sizes: A,B and C; and corresponding $nN \times nT$ configurations). Symbol “o” indicates the mean execution time and “*” corresponds to the maximum and minimum execution time.

somehow opposite behaviour of EULER and its higher sensitiveness suggest to accomplish a further in depth kernel-by-kernel analysis. A comparison of the different behaviour found in ACME and EULER in terms of the execution time for the two MG and EP kernels (memory- and CPU-intensive, respectively) is depicted in Fig. 8 for Dedicated Nodes setup. The plots for the MG kernel with classes B and C show that the speedup increases with monotone trend as more nodes are involved in the $nN \times nT$ configuration ($nT = 4, 8, 16$ and 32). And it is seen that this speedup is significantly greater in EULER, which can be explained because EULER nodes are more memory-bound than ACME's. The EP kernel in ACME exhibits also a rather small speedup when tasks are distributed over nodes. But on the contrary, EULER shows that the EP kernel runs slower when it is taken “out of the box” (that is, when the tasks are partially migrated from all being grouped in one node). It is visible in Fig. 8 that a big jump in execution time occurs as the EP kernel goes from

$1 \times nT$ to $2 \times nT$ (with $nT = 4$ and 8). It is noticed that the computed cases in EULER corresponding to $nT = 16$ and $nT = 32$ start at configurations 2×8 and 4×8 , respectively, thus it is not possible to have evidence of the "out of the box" effect in these cases (but it is plausible that the wavy pattern observed in these be similar to the wavy one observed for $nT = 8$ from the 2×4 configuration on). In summary, the rule derived for the EP kernel in EULER is that MPI task grouping makes sense as the execution time drops. And besides, the reversed behaviour seen in ACME for the EP kernel can be justified because of the much higher internode bandwidth, which compensate the "out of the box" effect observed in EULER.

4.4 Sensitivity to the clusters setup

The performance of kernels MG and EP is plotted in Fig. 9 corresponding to 8 MPI tasks and the four clusters setups analysed for the classes A, B and C. This plot is relevant because it provides four $nN \times nT$ configurations that start at 1×8 in both clusters, so the "out of the box" effect can be focused, if any. For both kernels, Dedicated Network and Dedicated Nodes setups provide very similar execution time, showing a monotone, quasi-linear drop with the number of nodes, which suggests the benefit (but small) of adding nodes to the computation. The plot for Dedicated Sockets is similar, but even it is seen a smaller drop of the execution time. This pattern is observed for the MG and EP kernels with all classes (it is noted that Dedicated Network and Dedicated Sockets setups are included only for ACME since EULER is a production cluster and only a portion of it was assigned to this research). Again, for the MG kernel in EULER under Dedicated Nodes, it is observed a higher speedup with the number of nodes, compared to ACME. In particular, for the execution of the EP kernel in EULER, it is visible the out-of-the-box effect under both setups: a large increase of the execution time when the EP kernel goes from 1×8 to a 2×4 configuration, followed by a saturation of the drop when additional nodes are included. A different conclusion is derived for EULER: while for the memory-bound MG kernel is beneficial to distribute the MPI tasks over so many nodes as possible (the smaller host memory of its sockets may explain the significant improvement), the CPU-bounded EP kernel demands to group them into one node to attain the best performance. In summary, the MG and EP kernels under Dedicated cores setup in ACME, points out to the dominance of a performance drop (but not monotone) when additional nodes are added to the computation (the execution time shows a pattern of either a moderate increase of up to 20%, or a slight drop in some cases). Comparison of the speedup obtained with the Dedicated Nodes and Dedicated Cores setups for the whole set of experiments conducted in ACME and EULER is given in Fig. 10. A comparison for ACME with 3 setups is provided in Fig. 11. The plotted boundary lines indicate the unused portion of the cluster due to the Slurm setup itself, which serves to suggest a general criterion about how much speedup is possible and which is the extra cost due to not using a portion

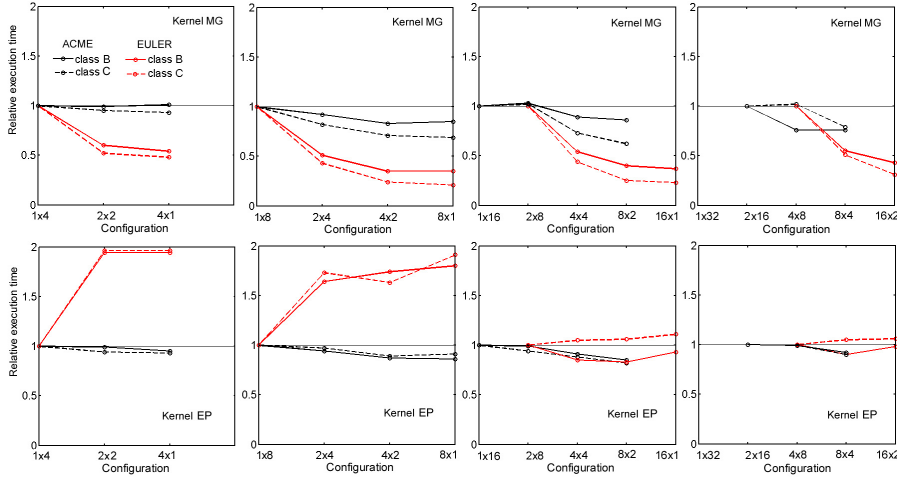


Fig. 8 Relative execution time for the Dedicated Nodes setup. NAS kernels MG (upper row) and EP (lower row) are shown for clusters ACME and EULER.

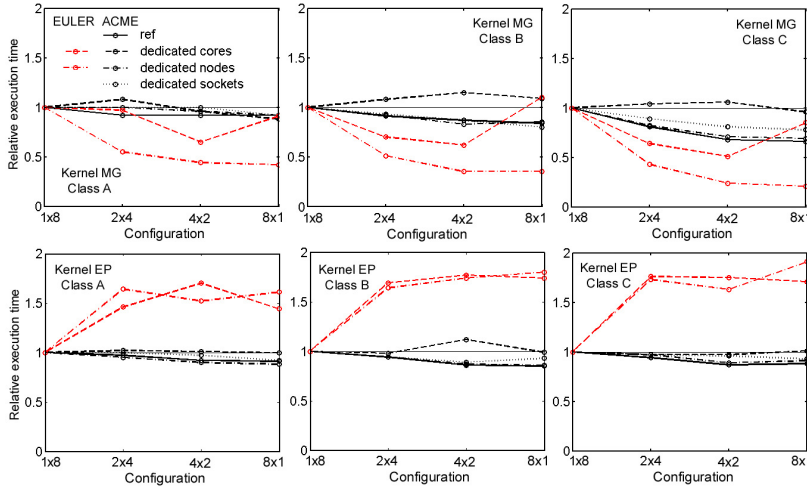


Fig. 9 Impact of cluster setup on the relative execution time for MG (memory-intensive) and EP (CPU-intensive) kernels with 8 MPI tasks in clusters.

of the machine (e.g.: say an ACME 2x4 configuration with Dedicated Nodes setup. This implies 4 tasks running on a socket of 8 cores. Since each node has 2 sockets, the occupation reads 4 of a total of $8+8=16$ cores, which means a 75% of unused cores).

The vertical scales in the plots relate speedup and % of unused cores (i.e. speedup of 2 corresponds to a 50% of unused cores; speedup of 4 to a 75% of unused cores; and so on). Fig. 10 shows that speedup values dispersion is greater in EULER. Speedup dispersion in ACME strongly increases at higher

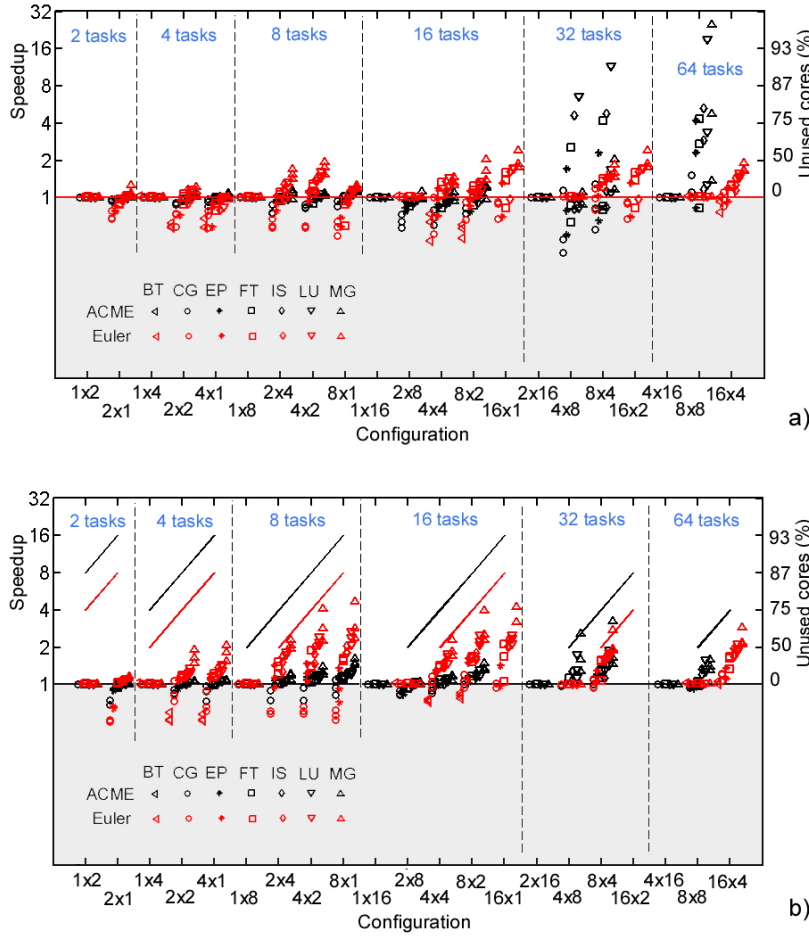


Fig. 10 Map of speedup VS. usage of computational resources for NAS under a) Dedicated Cores setup; and b) Dedicated Nodes setup (locus of % of unused cores is plotted for each number of MPI tasks. Line color corresponds to the symbols), depicted for clusters ACME and EULER.

degrees of parallelism (32 and 64 MPI tasks), mostly with the Dedicated Cores and Dedicated Sockets setup, as it is visible in Fig. 11.

These plots remark the importance of searching for a balance between significant speedups and not having too many unused cores. Obviously, it is a matter of settling a sweet point for users and cluster administrators. But under the Dedicated Nodes setup, it is seen that few points are over the boundary lines. A question that arises is how to cast a criterion to anticipate how the execution of a given experiment should be to speed it up as much as possible. Because only the Dedicated Cores setup has all its boundary lines collapsed into the 0%-unused cores situation (full occupation), it is interesting

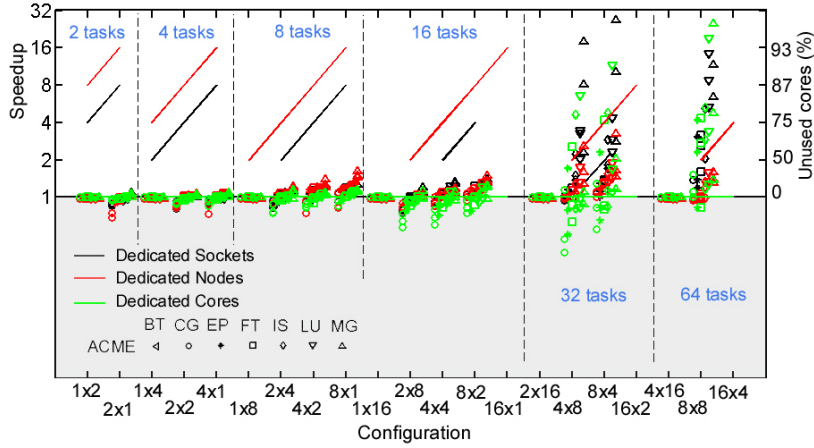


Fig. 11 Map of speedup VS. usage of computational resources for NAS under Dedicated Cores, Sockets and Nodes setups in cluster ACME.

to compare it with what happens using the other cluster setups. Let

$$S_{increment} = Speedup_i - Speedup_{DC}$$

be the speedup increment when the cluster setup is changed to $i=DN, DS$ (DN: Dedicated Nodes; DS: Dedicated Sockets). If

$$S_{increment} > 0$$

then the i -setup is potentially a better option for performance; but the context of unused cores in the cluster must be taken into account. For that, a speedup figure of merit (SFM) is introduced according to

$$SFM = \frac{Speedup_i}{Speedup_{DC}}$$

If the condition

$$SFM > \frac{total_{cores}}{used_{cores}}$$

is fulfilled by an executed experiment in the cluster, the attained performance improvement justifies the cluster operation with that number of idle cores. From a preliminary standpoint, this figure of merit may be useful to set a criterion to run production codes using an a priori building of a database out of code executions. That database would contain values of $Speedup_i$ and $Speedup_{DC}$ obtained under a variety of $nN \times nT$ configurations, problem sizes and short enough runtimes to be effective in the code characterization. Hence, this database would be a reduced "map of cases" to inferred executions behaviour and to guide throughput decisions.

5 Conclusions

The NAS Parallel Benchmarks have been executed in a systematic way on two clusters with rather different internode and intranode bandwidth properties, to identify dependencies between MPI tasks mapping and execution time speedup or resource occupation. The study comprises jobs up to 128 MPI tasks, justified by the limited strong-scaling properties of NAS kernels. A criterion for decision on which cluster setup is better to optimise the total cluster throughput is provided for NAS and its extension to production codes is suggested.

The findings can be related to two scenarios. When the clusters are configured to run parallel jobs in exclusivity, the results show that in most cases the execution time drops as the MPI tasks are distributed over the nodes (this agrees with previous investigations) and it seems efficient to distribute a given parallel job over cores located in different nodes. However, the other important scenario found in production HPC clusters corresponds to the need of sharing resources among set of jobs, such that the socket cores execute MPI tasks of different jobs. In this situation, a rather different behaviour is observed, much more sensitive to the type of cluster. In the state-of-the-art cluster ACME, many of the carried out experiments show a speedup when MPI tasks run on the fewest number of nodes. This is opposite to what is found in the cluster EULER, where execution time trends are more sensitive to the algorithm properties and part of the experiments points out to task distribution over nodes to shorten the execution time. This major result found in ACME feeds the discussion about possible computational efficiency benefits by tailoring live tasks migration and scheduling policies in modern clusters. In production clusters which share a significant load of serial jobs while running parallel jobs (a 7-year analysis of the executed tasks in EULER showed that more than half were serial), the question of to which extent serial tasks may act as perturbations to the execution of parallel jobs arises and deserves consideration to clarify the best live task migration policies within the context of optimizing clusters occupation. Other interesting aspect is how different the results would be in the case of hybrid MPI/OpenMP tasks. These open issues are part of the ongoing research within our group.

Acknowledgements This work was partially funded by the Spanish Ministry of Economy, Industry and Competitiveness project CODEC2 (TIN2015-63562-R) with European Regional Development Fund (ERDF) as well as carried out on computing facilities provided by the CYTED Network RICAP (517RT0529).

References

1. Bailey, D., et al.: The NAS Parallel Benchmarks. Tech. rep. (1994)
2. Bonnie++: www.coker.com.au/bonnie++
3. Chai, L., Gao, Q., Panda, D.K.: Understanding the impact of multi-core architecture in cluster computing: A case study with Intel dual-core system. In: Proc.- 7th IEEE Int. Symp. Cluster Computing and the Grid, CCGrid, pp. 471–478 (2007)

4. Chavarría-Miranda, D., Nieplocha, J., Tipparaju, V.: Topology-aware tile mapping for clusters of SMPs. *Proc. 3rd Conf. on Computing Frontiers* **2006**, 383 (2006)
5. Chi Zhang, Xin Yuan, A.S.: Processor Affinity and MPI Performance on SMP-CMP Clusters. In: *IEEE Int. Symp. Parallel & Distributed Processing, Workshops and PhD Forum*, pp. 1–8. Atlanta, USA (2010)
6. Intel Memory Latency Checker 3.1: www.intel.com/software/mlc
7. Jeannot, E., Mercier, G., Tessier, F.: Process Placement in Multicore Clusters: Algorithmic Issues and Practical Techniques. *IEEE Transactions on Parallel and Distributed Systems* **25**(4), 993–1002 (2014)
8. McCalpin, J.D.: Memory Bandwidth and Machine Balance in Current High Performance Computers. *IEEE TCCA Newsletter* (May), 19–25 (1995)
9. OSU Micro-Benchmarks: <http://mvapich.cse.ohio-state.edu/benchmarks>
10. Ribeiro, C.P., et al.: Evaluating CPU and Memory Affinity for Numerical Scientific Multithreaded Benchmarks on Multi-cores. *IJCSIS* **7**(1), 79–93 (2012)
11. Rodrigues, E.R., Madruga, F.L., Navaux, P.O.A., Panetta, J.: Multi-core aware process mapping and its impact on communication overhead of parallel applications. *Proc.-IEEE Symp. Computers and Comm.* pp. 811–817 (2009)
12. Shainer, G., Lui, P., Liu, T., Wilde, T., Layton, J.: The impact of inter-node latency versus intra-node latency on HPC applications. In: *Proc. IASTED Int. Conf. Parallel and Distributed Computing and Systems*, pp. 455–460 (2011)
13. Smith, B., Bode, B.: Performance effects of node mappings on the IBM BlueGene/L machine. *Euro-Par 2005 Parallel Processing* pp. 1005–1013 (2005)
14. Top 500: www.top500.org
15. Wu, X., Taylor, V.: Performance modeling of hybrid MPI/OpenMP scientific applications on large-scale multicore supercomputers. *J. Computer and System Sciences* **79**(8), 1256–1268 (2013)
16. Xingfu, W., Taylor, V.: Processor Partitioning: An Experimental Performance Analysis of Parallel Applications on SMP Clusters Systems. In: *19th Int. Conf. Parallel Distributed Computing and Systems (PDCS'07)*, pp. 13–18. CA, USA (2007)
17. Xingfu, W., Taylor, V.: Using Processor Partitioning to Evaluate the Performance of MPI, OpenMP and Hybrid Parallel Applications on Dual- and Quad-core Cray XT4 Systems. In: *Cray UG Proceedings (CUG 2009)*, pp. 4–7. Atlanta, USA (2009)