

Teaching HPC Systems and Parallel Programming with Small Scale Clusters of Embedded SoCs

Lluc Alvarez
Barcelona Supercomputing Center
Universitat Politècnica de Catalunya
lluc.alvarez@bsc.es

Eduard Ayguade
Barcelona Supercomputing Center
Universitat Politècnica de Catalunya
eduard.ayguade@bsc.es

Filippo Mantovani
Barcelona Supercomputing Center
filippo.mantovani@bsc.es

Abstract—In the last decades, the continuous proliferation of High-Performance Computing (HPC) systems and data centers has augmented the demand for expert HPC system designers, administrators and programmers. For this reason, most universities have introduced courses on HPC systems and parallel programming in their degrees. However, the laboratory assignments of these courses generally use clusters that are owned, managed and administrated by the university. This methodology has been shown effective to teach parallel programming, but using a remote cluster prevents the students from experimenting with the design, set up and administration of such systems.

This paper presents a methodology and framework to teach HPC systems and parallel programming using a small-scale cluster of embedded System-on-Chip (SoC) boards. These SoCs are very cheap, their processors are fundamentally very similar to the ones found in HPC, and they are ready to execute Linux out of the box, so they provide a great opportunity to be used in laboratory assignments for the students to experience with assembling a cluster, setting it up, and configuring all the software ecosystem. In addition, this paper shows that the small-scale cluster can be used as the evaluation platform for parallel programming assignments.

Index Terms—HPC systems, parallel programming, teaching

I. INTRODUCTION

The importance of High-Performance Computing (HPC) in our society has continuously increased over the years. In the early years, the very few existing HPC systems were based on vector processors specialised for scientific computations and they were only used by a small amount of experts; programmability and usability were not the key issues at that moment. The trend changed when supercomputers started to adopt "high-end" commodity technologies (e.g. general-purpose cores), which opened the door to a rich software ecosystem and, consequently, to many advantages in programming productivity. This was a key factor for the popularisation of HPC infrastructures, which spread throughout many research and industrial sectors. In the last years, the proliferation of HPC systems and data centers has gone even further with the emergence of mobile devices and cloud services. In the current scenario, the demand for expert HPC system designers, administrators and programmers is higher than ever, and will

This work is partially supported by the Spanish Government through Programa Severo Ochoa (SEV-2015-0493), by the Spanish Ministry of Science and Technology (contracts TIN2015-65316-P and FJCI-2016-30985), and by the Generalitat de Catalunya (contract 2017-SGR-1414).

likely continue growing to keep improving the performance and efficiency of HPC systems in the future.

In the last years, many universities have introduced courses on HPC systems and parallel programming in their degrees. Given the cost of modern HPC infrastructures, the laboratory assignments of most of these courses use clusters that are owned, managed and administrated by the university. This methodology is convenient to teach parallel programming, as the students only need to connect remotely to the cluster to do the programming work for the assignment. However, using a remote cluster prevents the students from experimenting with the design, set up and the administration of such systems.

Fortunately, with the advent of Systems-on-Chip (SoCs) for the embedded and the multimedia domains, today it is very easy and very cheap to build a small-scale cluster. Modern commercial SoCs for the embedded domain equip processors that are fundamentally very similar to the ones found in HPC systems, and are ready to execute Linux out of the box. So, these devices provide a great opportunity for the students to experience with assembling a cluster, setting it up, and configuring all the required software to have a fully operative small-scale HPC cluster.

This paper presents the methodology and framework that we propose for teaching HPC systems and parallel programming using a small-scale HPC cluster of embedded SoCs. This methodology has been successfully used to support teaching activities in Parallel Programming and Architectures (PAP), a third-year elective subject in the Bachelor Degree in Informatics Engineering at the Barcelona School of Informatics (FIB) of the Universitat Politècnica de Catalunya (UPC) - BarcelonaTech. After presenting the PAP course description and environment we first give an overview of the components of the small-scale cluster, which we name *Odroid cluster* after the Odroid-XU4 boards [1] that form it. Then the paper describes the methodology that we use in two laboratory assignments of the course. The first laboratory assignment consists on setting the Odroid cluster up and doing an evaluation of its main characteristics. The cluster setup consists on physically assembling the boards, configuring the network topology of the cluster, and installing all the software ecosystem typically found in HPC platforms. In the evaluation part the students discover the main characteristics of the Odroid-XU4 boards, they learn how the threads and

the processes of a parallel program are distributed among the processors and the nodes, and they experiment with the effects of heterogeneity. The second laboratory assignment consists on parallelizing an application implementing the heat diffusion algorithm with MPI [2] and OpenMP [3] and evaluating it on the Odroid cluster. The complete framework presented in this paper greatly facilitates the learning of the design, the set up and the software ecosystem of HPC systems, as well as being a very appealing platform for the evaluation of parallel programming assignments.

The rest of this paper is organized as follows: Section II explains the course and its methodology. Section III gives an overview of the Odroid cluster and its components. Section IV describes the work to be done by the students to evaluate the Odroid cluster and to understand its main characteristics. Section V then shows how we use the Odroid cluster as a platform for a parallel programming assignment. Section VII remarks the main conclusions of this work. Finally, section VIII presents as an annex the step-by-step process that is followed by the students in the laboratory assignment to set up the cluster.

II. CONTEXT, COURSE DESCRIPTION AND METHODOLOGY

Parallel Programming and Architectures (PAP) is a third-year (sixth term) optional subject in the Bachelor Degree in Informatics Engineering at the Barcelona School of Informatics (FIB) of the Universitat Politècnica de Catalunya (UPC) - BarcelonaTech. The subject comes after Parallelism (PAR [4]), a core subject in the Bachelor Degree that covers the fundamental aspects of parallelism, parallel programming with OpenMP and shared-memory multiprocessor architectures. PAP extends the concepts and methodologies introduced in PAR, by focussing on the low-level aspects of implementing a programming model such as OpenMP, making use of low-level threading (*Pthreads*); the subject also covers distributed-memory cluster architectures and how to program them using MPI. Another elective course, *Graphical Units and Accelerators* (TGA) explores the use of accelerators, with an emphasis on GPUs, to exploit data-level parallelism. PAR, PAP and TGA are complemented by a compulsory course in the Computer Engineering specialisation, *Multiprocessor Architectures*, in which the architecture of (mainly shared-memory) multiprocessor architectures is covered in detail. Another elective subject in the same specialisation, *Architecture-aware Programming* (PCA), mainly covers programming techniques for reducing the execution time of sequential applications, including SIMD vectorisation and FPGA acceleration.

The course is done in a 15-week term (one semester), with 4 contact hours per week: 2 hours dedicated to theory/problems and 2 hours dedicated to laboratory. Students are expected to invest about 5-6 additional hours per week to do homework and personal study (over these 15 weeks). Thus, the total effort

devoted to the subject is 6 ECTS credits¹.

The content of the course is divided in three main blocks. The first block has the objective of opening the black box behind the compilation and execution of OpenMP programs; it covers the internals of runtime systems for shared-memory programming, focusing on the most relevant aspects of thread management, work generation and execution, and synchronisation; in a very practical way, students explore different alternatives for implementing a minimal OpenMP-compliant runtime library, using *Pthreads*, providing support for both the work-sharing and tasking execution models. This block takes 4 theory/problems sessions (mainly covering low-level *Pthreads* programming) and 6 laboratory sessions (individual work). At the end of the laboratory sessions for this block, a session is devoted to share experiences and learnings. Additional details about this block and the laboratory assignment can be found elsewhere [5].

The second block has the objective of understanding the scaling path to parallel machines with large number of processors, beyond the single-node shared-memory architectures students are familiar with; it covers the main hardware components of such systems (processors, accelerators, memories and their interconnection). This block takes 4 theory/problems sessions devoted to analyse in detail how the ratio FLOPs/Byte evolves in the scaling path (i.e. number of potential floating-point operations per byte of data from accessed from/to memory/interconnect). The roofline model [6], plotting floating-point performance as a function of the compute units peak performance, data access peak bandwidth and arithmetic intensity, is used to understand this evolution and its implications on data sharing in parallel program; the evolution of the energy efficiency (Flops/Watt) is also covered. The block also takes 3 laboratory sessions in which students 1) physically assemble a small cluster based on Odroid boards [1], 2) set up the Ethernet network and the Network File System (NFS), 3) install and configure all the software required to execute MPI and OpenMP parallel programs, and 4) evaluate the cluster using some simple benchmarks. This laboratory work is complemented with a guided learning assignment in which groups of students propose the design of a real HPC system, based on commodity components, with certain performance/power trade-offs and economic budget; this is an excellent opportunity for them to take a look at real components and include cost as one of the important trade-offs in HPC system design. The proposed designs are presented, discussed and ranked in a session with the idea of sharing the criteria used by each group of students.

The last block in the course has the objective of studying the basics of parallel programming for distributed-memory architectures, using MPI; it covers the aspects related with the creation of processes/groups and the different data com-

¹The European Credit Transfer System (ECTS) is a unit of appraisal of the academic activity of the student. It takes into account student attendance at lectures, the time of personal study, exercises, labs and assignments, together with the time needed to do examinations. One ECTS credit is equivalent to 25-30 hours of student work.

munication strategies and the trade-offs. This block has a duration of 3 theory/problems and 3 laboratory sessions. In the laboratory students develop a hybrid MPI/OpenMP implementation of the classical heat diffusion problem, evaluating its performance on the Odroid cluster that they have already assembled. Although the laboratory assignment in this block could be done in production cluster available at the department/university, we preferred to continue in the Odroid cluster. During the 6 sessions students are expected to administrate the Odroid cluster. In particular, they are asked to deal with all the potential problems that they may encounter, to write scripts to automatise setup and evaluation processes, and to install libraries and tools (editors, debuggers, etc) to have a productive programming environment.

The rest of the paper focus on the laboratory activities related with these two last blocks in PAP.

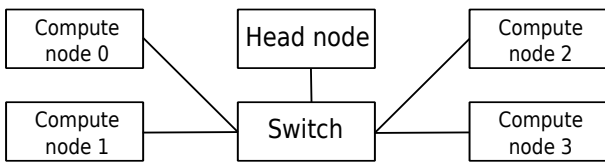


Fig. 1: Scheme of the Odroid cluster.

III. ODROID CLUSTER OVERVIEW AND COMPONENTS

This section provides an overview of the Odroid cluster that is used in the laboratory assignments. Figure 1 illustrates the main components of the Odroid cluster. It consists of one *head node* and four *compute nodes* connected through a switch. The *head node* acts as the gateway for accessing the cluster and is also in charge of providing Internet connection to the whole cluster and to host the DHCP server. The hardware components required to assemble the cluster are listed below:

- *head node*: personal computer (PC) with 2 network interfaces.
- *compute nodes*: 4 Odroid-XU4 boards, each with a eMMC card with a pre-installed Ubuntu 16.04, as shown in Figure 2.
- 1 8-port Gigabit Ethernet desktop switch.
- 1 power supply for the Odroid-XU4 boards and switch.

The Odroid-XU4 boards are based on the Samsung Exynos5 Octa chip [7], a low-power heterogeneous multicore processor based on the ARM big.LITTLE architecture [8]. Each processor consists of 4 Cortex-A15 out-of-order cores running at 1.9 GHz and 4 Cortex-A7 in-order cores running at 1.3 GHz. In addition the board includes a LPDDR3 RAM chip of 2GB as main memory and ports for Gigabit Ethernet, eMMC 5.0 and μ SD Flash storage, USBs and HDMI display. The board also comes with an active cooling fan mounted on top of the socket (not shown in Figure 2).

The *head node* and the 4 *compute nodes* are connected to the 8-port Gigabit Ethernet switch. A picture of the assembled cluster is shown in Figure 3. As shown in the picture, we provide a methacrylate plate that has the switch and the

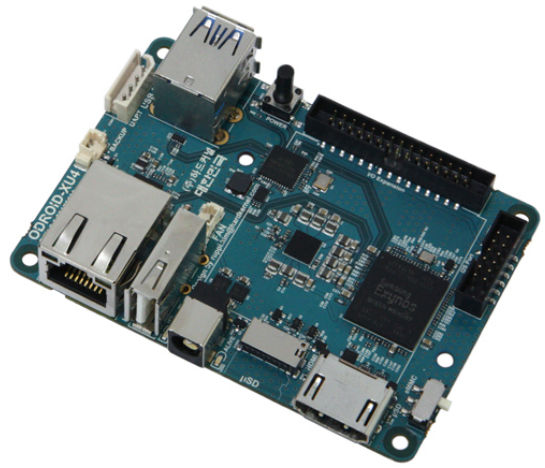


Fig. 2: Odroid-XU4 board.

power supply already attached, and also some free space for the students to stack the Odroid boards vertically. This way the cluster is much more compact and only a single plug is required for the whole cluster.

Section VIII is an annex that provides detailed information about the required steps to set the Odroid cluster up.

IV. ODROID CLUSTER EVALUATION

A. Cluster Characteristics

In order to understand the main characteristics of the Odroid boards, we ask the students to use the commands `lscpu` to obtain general information about the node and `"lscpu -e"` to get per-core information. We also encourage them to inspect the file `/cpu/cpuinfo` to get additional details about the cores. In addition, the students are asked to use the `lstopo` command (included in the `hwloc` package) to get a graphical representation of all the components of a compute node.

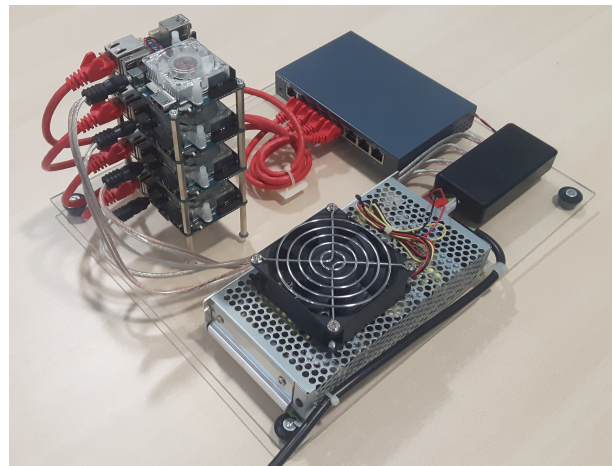


Fig. 3: Picture of the Odroid cluster.

B. Thread and Process Distribution

To find out how the threads and the processes are distributed across the Odroid cluster, we ask the students to experiment with a "Hello world!" benchmark programmed in hybrid MPI/OpenMP. The benchmark simply prints, for each thread, its process ID, its thread ID and the name of the MPI host where it is executed.

First the students experience with the MPI process distribution. To do so they disable the OpenMP parallelization at compile time and execute the hello world benchmark with the command `mpirun.mpich -np P -machinefile machines ./hello_world`. Note that the `-np` option specifies the number of MPI processes (P), and the `-machinefile` option specifies the file that describes the names of the nodes to be used in the execution and the number of processes per node. We provide two machine files, shown in Figure 4, one that specifies 8 processes per node and one that specifies 1 process per node. The students execute with 1 to 32 processes using both machine files and report how the processes are distributed across the nodes.

odroid-0:8	odroid-0:1
odroid-1:8	odroid-1:1
odroid-2:8	odroid-2:1
odroid-3:8	odroid-3:1

(a) 8 processes per node (b) 1 process per node

Fig. 4: Machine files for MPI programs.

Then the students experience with the distribution of MPI processes and OpenMP threads. To do so they compile the benchmark with support for OpenMP and execute it with the command `mpirun.mpich -np P -machinefile machines -genv OMP_NUM_THREADS T ./hello_world`. In addition to the number of MPI processes and the machine file, the `-genv` option specifies the number of OpenMP threads (T) that each MPI process spawns. The students test both machine files in executions with 1 to 4 processes and 1 to 8 threads and observe how these are distributed in the cluster.

C. Heterogeneity

One of the main characteristics of the Odroid boards is the heterogeneity of the cores, as explained in Section III. We provide the students with a very simple compute-intensive kernel that calculates the number Pi, which allows to clearly understand and experience the heterogeneity of the ARM big.LITTLE architecture.

Listing 1 shows the code of the Pi benchmark. The first part (lines 2 to 6) initializes the variables and calculates the start and end of the iteration space (for the loop in line 9) assigned to each MPI process. The second part (lines 8 to 12) is the main kernel, which is a compute-intensive loop that can be parallelized using the OpenMP `parallel for` construct with a reduction. The last part (lines 14 to 16) calculates the partial results (`local_pi`) in each MPI process and reduces

```

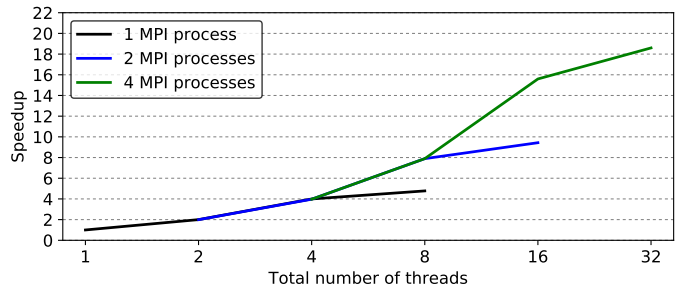
1 void pi(int num_steps) {
2     chunk_size = num_steps / num_procs;
3     start = proc_id * chunk_size;
4     end = (proc_id + 1) * chunk_size;
5     h = 1.0 / (double) num_steps;
6     sum = 0.0;
7
8     #pragma omp parallel for reduction(+: sum)
9     for (i = start; i < end; ++i) {
10         x = h * ((double)i - 0.5);
11         sum += 4.0 / (1.0 + x*x);
12     }
13
14     local_pi = h * sum;
15     MPI_Reduce(&local_pi, &global_pi, 1,
16             MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
17 }

```

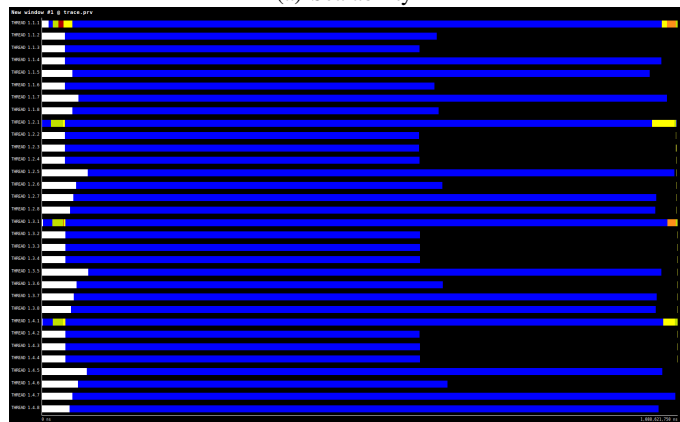
Listing 1: Pi source code.

them in a global variable (`global_pi`) to obtain the final result.

We ask the students to compile and execute the Pi benchmark with 1.28G steps and different number of MPI processes (1 to 4) and OpenMP threads (1 to 8). Note that the OpenMP `parallel for` construct in the source code does not specify any schedule for the iterations of the loop, so the static one is used by default (one contiguous block of $(end - start) \div T$ iterations per thread, being T the number of OpenMP threads per MPI process. Figure 5a shows the scalability of the bench-



(a) Scalability



(b) Execution trace

Fig. 5: Scalability and execution timeline of the Pi benchmark with a static OpenMP schedule.

mark with the static scheduler. The figure shows the speedup achieved by augmenting the total number of threads. The total number of threads are derived from the executions with 1 to 4 MPI processes and 1 to 8 OpenMP threads per process. The MPI processes are distributed across the nodes. The results show perfect scalability when up to 4 threads per process are used. This happens because the benchmark is compute intensive and the threads are scheduled on the fast cores by default. The parallel efficiency decreases when 8 threads per process are used, achieving speedups of 4.85x, 9.42x and 18.32x with 1, 2 and 4 MPI processes, respectively. These performance degradations are caused by the heterogeneity of the ARM big.LITTLE architecture. As shown in the execution timeline (obtained with *Extrae* and visualized with *Paraver*) in Figure 5b, the execution phases (in blue) of the fast cores end significantly earlier than the ones of the slow cores, so then the fast cores have to wait (in black) for the slow cores to finish. This execution imbalance (not work imbalance) is due to the different computing power of the two kinds of cores in the processor, preventing the Pi benchmark from scaling further in executions with more than 4 threads per process.

To mitigate the effects of heterogeneity we ask the students to try to improve the performance of the Pi benchmark by using a dynamic schedule in OpenMP with different chunk sizes (number of iterations per chunk dynamically assigned to a thread). Figure 6a shows the scalability of the Pi benchmark

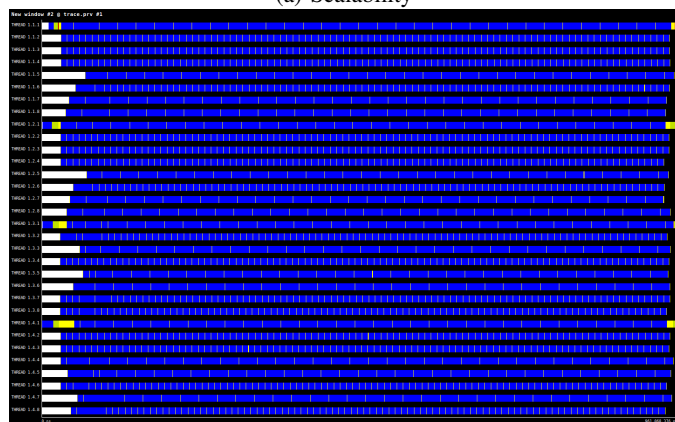
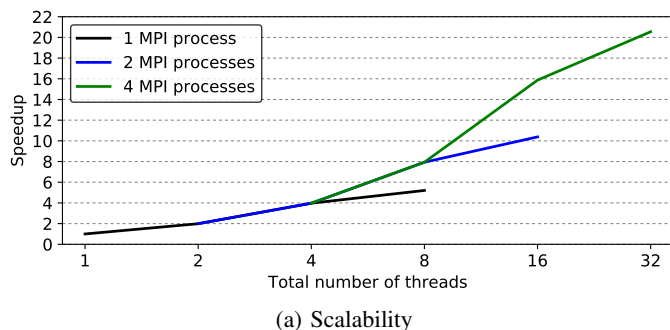


Fig. 6: Scalability and execution trace of the Pi benchmark with a dynamic OpenMP scheduler and chunk size 512.

with a dynamic scheduler and a chunk size of 512, which is the best granularity for this experiment. As in the case of the static scheduler, the results show perfect scalability with up to 4 threads per process. When 8 threads per process are used the scalability is slightly better than the one obtained with the static scheduler, achieving respective speedups of 5.26x, 10.34x and 20.50x with 1, 2 and 4 MPI processes. The execution trace in Figure 6b clearly shows how the dynamic scheduler perfectly distributes the work among the threads, so the time the threads spend waiting in the barrier is negligible. However, achieving perfect scalability is impossible because of two important aspects of the architecture. One is the reduced compute capabilities of the slow cores compared to the fast ones, and the second one is the DVFS controller that, in order not to exceed the power and temperature caps, lowers the frequency and the voltage of the cores when they are all active at the same time.

V. PARALLELIZATION OF THE HEAT DIFFUSION PROGRAM

In the last assignment of the course the students have to parallelize a sequential code that simulates heat diffusion in a solid body using the Jacobi solver for the heat equation. The program is executed with a configuration file that specifies the maximum number of simulation steps, the size of the bi-dimensional solid body and the number of heat sources, with their individual position, size and temperature. Two configuration files are provided, one for programming and debugging and one for evaluating the parallelization. The program reports some performance measurements (execution time, floating point operations, the residual and the number of simulations steps performed) and an image file with a gradient from red (hot) to dark blue (cold) for the final solution. Figure 7 shows the output of the program for a solid with two heat sources, one in the upper left corner and one at the bottom center.

An skeleton of the code for the heat diffusion benchmark is shown in Listing 2. Note that the code contains comments to guide the parallelization strategy, that is explained in the

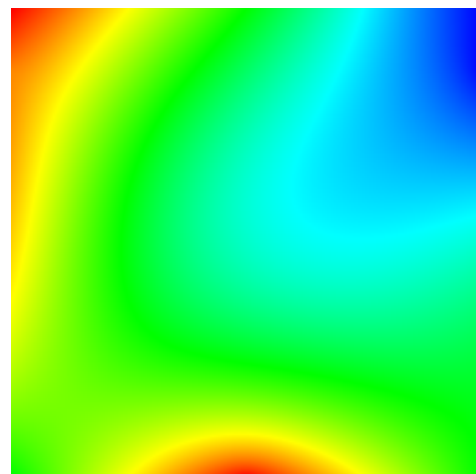


Fig. 7: Image of the temperature of the solid body.


```

1 void heat() {
2   Read configuration file
3   Allocate memory
4   Initialize matrices and heat sources
5
6   // MPI-1: Broadcast input parameters
7   // MPI-2: Distribute matrix
8
9   while(1) {
10
11     // MPI-3: Exchange halos
12
13     // OPENMP-1: Parallelize loop
14     for(i=1; i<N-1; i++) {
15       for(j=1; j<N-1; j++) {
16         utmp[i][j] = 0.20 * (
17           u[i][j] + // center
18           u[i][j-1] + // left
19           u[i][j+1] + // right
20           u[i-1][j] + // top
21           u[i+1][j]); // bottom
22
23         diff = utmp[i][j] - u[i][j];
24         residual += diff * diff;
25       }
26     }
27
28     // MPI-3: Exchange halos
29
30     aux = u;
31     u = utmp;
32     utmp = aux;
33
34     // MPI-4: Communicate residual
35
36     iter++;
37     if(residual < R && iter == MAX_ITERS)
38       break;
39   }
40
41   // MPI-2: Distribute matrix
42 }

```

Listing 2: Heat diffusion source code.

two next subsections. In the initialization (lines 2, 3 and 4), the program reads the input configuration file to establish the input parameters, including the size of the matrices (N), the threshold value for convergence (R), and the maximum number of iterations. Then the programs allocates the memory for two matrices (u and $utmp$) and initializes them according to the heat sources. Then the code enters the main loop (lines 9 to 39), that has three main parts. The first part (lines 14 to 26) computes one step of the heat diffusion simulation. The computation is a 5-point 2D stencil that uses two matrices, one as input and one as output. The second part (lines 30 to 32) swaps the matrices so that the output of the previous stencil becomes the input of the next stencil in the following iteration. The third part (lines 36 to 38) checks if the solution has converged or the maximum amount of iterations has been reached.

A. MPI Parallelization Strategy

The heat diffusion code we provide to the students with contains a commented skeleton for the MPI and the OpenMP

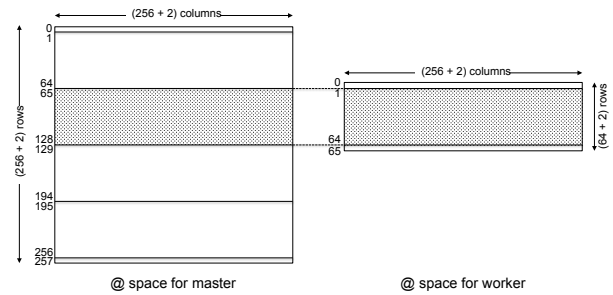


Fig. 8: Global and local storage for the matrix of the heat diffusion algorithm with different processes.

parallelization, as shown in Listing 2. The students first parallelize the application with MPI following a series of steps.

The first step, labeled as MPI-1 in line 6, consists on exchanging the information regarding the parameters of the execution between the processes. On one hand, the master first reads the configuration file, allocates memory and sends the execution parameters to the workers; then it computes the heat equation on the whole 2D space, and finally it reports the performance metrics and generates the output file. On the other hand, each worker receives from the master the information required to solve the heat equation, allocates memory, performs the computation on the whole 2D space and finishes. Note that this version does not benefit from the parallel execution since workers replicate the work done by the master.

The second step, labeled as MPI-2 in lines 7 and 41, consists on distributing the matrices among the processes so that the master and the workers solve the equation for a subset of consecutive rows. Figure 8 shows an example for a 256x256 matrix (plus the halos) distributed among 4 processes. After computing its part of the matrix, the workers return the part of the matrix they have computed to the master in order to re-construct the complete 2D data space. This version of the code does not generate a correct result because the processes do not communicate the boundaries to the neighbor processes at each simulation step.

The third step, labeled as MPI-3 in line 11 and 28, consists on adding the necessary communication so that, at each simulation step, the boundaries are exchanged between consecutive processors. The students have the freedom to use the MPI communication routines they find more appropriate. This version of the code still generates an incorrect solution because the total number of simulation steps done in each process is controlled by the local residual of each process and the maximum number of simulation steps specified in the configuration file, instead of computing the residual and checking the convergence globally.

The fourth and last step, labeled as MPI-4 in line 34, consists on adding the necessary communication to control the number of simulation steps of all the processes by computing a global residual value at every simulation step. The students again can use the communication strategy that they find most appropriate. This version of the code generates exactly the

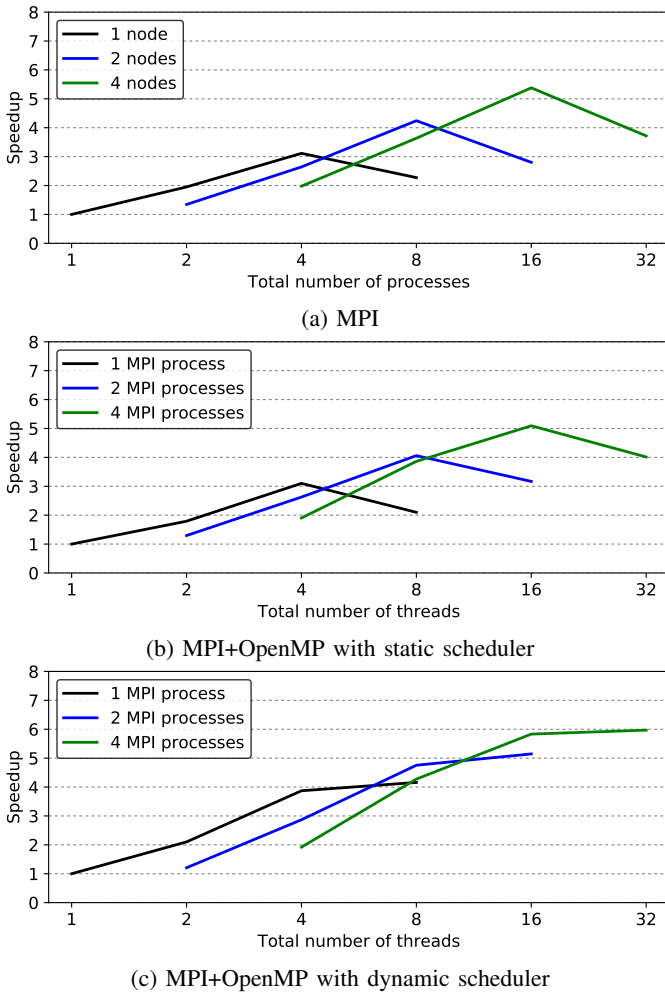


Fig. 9: Speedup of the different parallel implementations of the heat diffusion algorithm.

same result as the one generated by the sequential code.

B. Hybrid MPI/OpenMP Parallelization Strategy

Once the students have successfully parallelized the heat diffusion algorithm with MPI they proceed with the hybrid MPI/OpenMP implementation. Adding the OpenMP parallelization on top of the MPI implementation is quite trivial, as it only requires introducing an `OpenMP parallel for reduction(+:residual)` construct in the stencil kernel (labeled as `OPENMP-1` in line 13). We also encourage the students to try different OpenMP schedulers to mitigate the effects of the heterogeneity, as they have learned in the previous laboratory assignment.

C. Evaluation on the Odroid Cluster

After successfully parallelizing the heat diffusion program the students proceed with its performance evaluation on the Odroid cluster.

Figure 9a shows the scalability of the MPI parallel implementation of the heat diffusion benchmark. The figure shows the speedup achieved by augmenting the total number of

processes, which are derived from executions with 1 to 8 MPI processes per node and 1 to 4 nodes. It can be observed that, when up to 4 processes are used in each node, the program scales up to 3.06x, 4.17x and 5.41x with 1, 2 and 4 nodes, respectively. In addition, the performance drops with 8 processes per node, reducing the speedups to 2.16x, 2.89x and 3.81x with 1, 2 and 4 nodes, respectively. This exact trend is repeated in the hybrid MPI/OpenMP implementation with a static scheduler, as shown in Figure 9b. As in the case of the Pi benchmark, the performance degradations when using 8 threads or processes per node are caused by the heterogeneity of the cores and the inability of the parallel implementation to dynamically balance the load in such scenario. Figures 9a and 9b also reflect the higher cost of communicating data across the nodes rather than inside the nodes. It can be observed that, in the MPI version executing with 4 total processes, using 1 node for all the processes is 54% faster than spreading them across 4 nodes.

The scalability of the hybrid MPI/OpenMP implementation with a dynamic scheduler is shown in Figure 9c. It can be observed that, when using 8 threads per node, the performance does not drop as in the previous two versions of the code. However, the performance gains are very low compared to the executions with 4 threads per node and the same number of nodes. Another important difference is that the scalability obtained with up to 4 threads per node is slightly higher than in the other two implementations. With 4 threads per node the results show that the dynamic scheduler is 18%, 10% and 8% faster than the static one in executions with 1, 2 and 4 nodes, respectively.

VI. PREVIOUS AND RELATED WORK

The widespread availability and low cost of single-board computers based on SoC have provided to educators the possibility of building small clusters and explore different ways to use them in their courses related to parallel and distributed computing. Several "microclusters" based on various Raspberry Pi and Odroid models, Nvidia's Jetson TK1 or Adapteva's Parallella are presented in [9]. Authors in that publication also present the various strategies followed by them regarding the use of their microclusters and summarise some earlier examples of microclusters that have been inspirational for the recent proposals based on inexpensive single-board computers and.

In this paragraph we comment on some of the microclusters presented in [9] with the aim of observing the similarities and differences with the proposal in this paper. David Toth introduced in 2014 the first version of the Half Shoebox Cluster (HSC) [10], which equipped two compute nodes based on ARM Cortex-A7 dual core processors, with the purpose of teaching Pthreads, OpenMP and MPI. The HSC has been continuously evolving, and the six subsequent versions of the HSC were built with various Odroid SoCs such as the U3, the C1, the XU3-Lite, the C1+, the XU4, and the C2. The XU3-Lite and the XU4 have 8-core ARM CPUs, while the rest of Odroid boards have 4-core CPUs. Rosie [11], constructed in

2014 by Libby Shoop, consisted of 6 NVIDIA TK1 single board computers. each equipped with a quad-core Cortex-A15 processor and an integrated Kepler GPU with 192 cores. The cluster has been used to teach heterogeneous computing techniques with OpenMP, MPI and CUDA. Finally, Student-Parallella [12], with 4 Parallella nodes, each with a Zynq 7000-series dual core ARM Cortex-A9 and an Epiphany coprocessor with 16 cores, was constructed by Suzanne Matthews and covering the native Epiphany programming model, Pthreads, OpenMP and MPI in her parallel computing elective course.

Most of the previously mentioned microclusters were designed with the idea of having an alternative to 1) non-dedicated networks of workstations in laboratory classrooms, 2) the relatively expensive to build and maintain high-performance departmental/university clusters or 3) the "far-away" cloud systems such as Amazon's EC2. In addition, seeing the cluster in action, with all components and cables, encourages students interest for parallel and distributed concepts. Usually microclusters are given to the students pre-assembled and configured, ready to be used to learn parallel programming using Pthreads, OpenMP, MPI and/or CUDA.

Based on our experience at the Barcelona Supercomputing Center (BSC-CNS), hosting the Marenstrum cluster-in-a-chapel supercomputer, observing how well trained HPC system administrators were offered good positions to setup and administrate datacenters in companies and research institutions, we decided to design a module, practical in nature, to teach these skills. The first laboratory assignment presented in this paper was created making emphasis on the cluster setup process, HPC software ecosystem configuration and initial performance testing, trying to mimic as much as possible the real systems found in supercomputing centers. Contrary to the use of microclusters, the use of traditional HPC clusters or cloud systems hides all these aspects from students. Once build and tested by each group of students, the microcluster is programmed using hybrid MPI/OpenMP in the proposed second assignment.

Our first Odroid-XU3 microcluster was build during the spring semester in 2015 by a couple of undergraduate students doing an optional laboratory assignment; the initiative was followed by the current Odroid-XU4 microcluster that has been used since then. In fact, using the microcluster to learn and practice MPI programming (the second assignment) was a proposal that originated from our students after building, configuring and testing the microcluster during the 2016 spring semester.

There are organizations like XSEDE [13] and the Blue Waters program [14] that provide parallel computing resources for educational purposes. However, getting set up on these systems is a non-trivial process and requires the faculty member to be well-organized to submit a small grant proposal at the right time, get student accounts set up, and get comfortable with the systems themselves. For faculty new to parallel computing, these challenges can be a high barrier to being able to teach parallel computing.

VII. CONCLUSIONS AND EVOLUTION

The continuous expansion of HPC systems and data centers has come together with an increasing demand for expert HPC system designers, administrators and programmers. To fulfill this demand, most university degrees have introduced courses on parallel programming and HPC systems in recent years. However, very often the laboratory assignments of these courses only focus on the parallel programming part, and the students never experiment with the design, set up and administration of HPC systems.

This paper presents a methodology and framework to use small-scale clusters of single-board embedded SoCs to teach HPC systems and parallel programming. In contrast to the traditional methodology for teaching parallel programming with remote clusters managed by the university, using small-scale clusters allows the students to experience with assembling a cluster, setting it up, configuring all the software ecosystem, and administrating it during the duration of the course. In this paper we show that these SoCs have very appealing characteristics for being used in laboratory classes, given their low cost, their ability to execute the same software stack as HPC systems, and the similarity between their processors and the ones used in HPC. Moreover, we show that these small-scale clusters are also a very attractive platform to experience with relevant aspects of today HPC systems such as heterogeneity (different kinds of cores), variation of frequency with number of active cores, or cost of data communication.

A very positive adoption of the two assignments in the elective Parallel Architectures and Programming (PAP) course has been observed. Between 15 and 20 students yearly follow the course and build 5 independent clusters working in groups of 4 students; although the number of students per group may seem large, the curiosity and surprises found while doing the two assignments motivates interesting discussions among them.

VIII. ANNEX: ODROID CLUSTER SETUP

This section describes the required steps that to set the Odroid cluster up. These steps are followed by the students in the laboratory assignment of the second block of the course.

A. Assembling the Odroid Cluster

The first step is to identify all the components that we provide, which are listed in Section III (with all cables, screws and separators required). The students then physically assemble the Odroid boards, stacking them vertically using the separators.

B. Head Node Setup

The PC that is used as head node has two network interfaces and runs a Ubuntu Desktop 16.04. The primary functions of the head node are to share the Internet connection with the compute nodes and to act as the DHCP server.

To set the head node up the PC has to be connected to Internet using one of the network interfaces and to the Gigabit Ethernet switch using the second network interface. To share

the Internet connection with the compute nodes the students use the Linux Network Manager, which allows to easily share the Internet connection using a GUI. To do so, the students need to identify the network interface that is connected to the switch, edit its properties, and select the option “Method: Shared to other computers” in the “IPv4 Settings” tab.

To check that the head node is properly connected to Internet and to the switch, we ask the students to use and explain the output of the `ifconfig` command. If the connections are properly configured, the output shows that the head node is connected to two networks, one using the network interface that is connected to Internet (`enp0s25`, with IP 192.168.60.XXX by default), and one using the interface that is connected to the switch (`enp7s4`, with IP 10.42.0.1 by default).

C. Compute Nodes Setup

Each compute node boots its own operating system image from the eMMC card and is assigned an IP address by the DHCP server in the head node. The following steps need to be followed to appropriately configure the compute nodes as part of the cluster.

The first step to be done is to install the 4 eMMC cards in the 4 Odroid boards and to connect the 4 boards to the switch using 4 Ethernet cables. The boards include a switch to choose the boot media (μ SD or eMMC), so the 4 switches have to be positioned to boot from eMMC card. The eMMC cards that are provided to the students already contain a pre-installed Ubuntu 16.04 operating system.

The second step is to boot the compute nodes and to give them a unique host name. It is very important to turn on the compute nodes one by one because, since they all run the same operating system image, they all have the same host name, so they cause host name conflicts in the network if all of them try to boot at the same time. To boot a compute node and change its host name the students turn on the board and wait for it to boot. The compute nodes are not connected to any display nor peripheral, so all the interaction with them has to be done via SSH from the head node. To check if the compute node has booted and is connected to the network, the command `nmap -sn 10.42.0.0/24` has to be used from the head node. If the compute node is up, the output of this command shows a line similar to this: `Nmap scan report for 10.42.0.230; Host is up (0.0012s latency)`. This means that the compute node is up and is connected to the network with the IP address 10.42.0.230, so the students can connect from the head node to the compute node using SSH to its IP address (`ssh odroid@10.42.0.230`). Once connected to the compute node, changing its host name can be done by simply editing the file `/etc/hostname` so that it contains the desired name. After changing the host name the board needs to be rebooted and, after checking that the compute node has booted correctly and is visible in the network (using again the `nmap` command from the head node), the student can repeat this step for the rest of compute nodes.

We encourage the students to use a naming convention for the compute nodes that facilitates their identification in the physical rack. For instance, we propose to name the compute nodes as `odroid-0`, `odroid-1`, `odroid-2` and `odroid-3`, being `odroid-0` the board at the bottom and `odroid-3` the board at the top of the rack.

D. Trusted SSH Connections

SSH keys are a well-known way to identify trusted computers in a network. We use SSH keys in the Odroid cluster to allow any compute node to access any other one without requiring any password. In order to do so, the students have to follow the next steps for each compute node.

The first step is to generate a private authentication key pair (a public and a private key) using the command `ssh-keygen -t rsa`. The generated public and private keys are by default located in `.ssh/id_rsa.pub` and `.ssh/id_rsa` in the home directory, respectively.

The second step is to add the public key to the list of authorized keys for the same compute node. To do so the students simply need to execute the command `cat .ssh/id_rsa.pub >> .ssh/authorized_keys`.

The third and last step is to transfer the public key generated for the compute node to the rest of compute nodes. This is done with the command `ssh-copy-id odroid@odroid-X`, being `odroid-X` the compute node to which the public key is transferred. So, to transfer the public key of a compute node (i.e., `odroid-0`) to the rest of compute nodes, the command has to be executed three times (i.e., to `odroid-1`, to `odroid-2` and to `odroid-3`).

These previous three steps have to be repeated in all the compute nodes. At the end of the process each compute node should have the public keys of all the compute nodes, and any compute node should be able to access any other one without entering any password. We ask the students to make sure the whole process worked by trying to access different compute nodes from each other via SSH.

E. NFS Setup

Network File System (NFS) is a distributed file system protocol that allows different nodes to share files over a network. The protocol requires one node to act as a NFS server, while the rest of nodes act as clients. In this section we assume the compute node `odroid-0` is the NFS server, while `odroid-1`, `odroid-2` and `odroid-3` are the NFS clients. In order to set up a shared directory between the compute nodes the students follow a series of steps.

The first step is to install the NFS packages and to create the directory that is going to be shared across the nodes. This step is done in all the compute nodes. To install the NFS packages the students just need to use the command `apt install nfs-common`, and then they create the directory that will be shared via NFS, i.e. `/sharedDir`.

The second step is to configure one of the compute nodes to be the NFS server. To do that, the students first have to select one of the compute nodes to be the NFS server, `odroid-0`

for example. Then they have to install the NFS server packages in the selected node with the command `apt install nfs-kernel-server`. Finally they have to export the NFS directory of the NFS server node to the rest of nodes by editing the file `/etc/exports` on `odroid-0` so that it contains the following: `/sharedDir *(rw, sync)`. Once this is done the NFS server has to be restarted with the command `sudo service nfs-kernel-server restart`.

The third step is to configure the NFS directory in the rest of compute nodes. Assuming the compute node `odroid-0` is the NFS server and the rest of nodes are the NFS clients, the students can mount the directory `/sharedDir` of the compute node `odroid-0` on a local directory of the compute nodes `odroid-1`, `odroid-2` and `odroid-3` using the command `mount -t nfs odroid-0:/sharedDir /sharedDir`. However, we encourage the students to automatize this so that it gets mounted at boot time. This can be done by modifying the file `/etc/fstab` of the NFS client nodes (`odroid-1`, `odroid-2` and `odroid-3`) so that it contains the line `odroid-0:/sharedDir /sharedDir nfs`. Note that the file `/etc/fstab` is formed by columns. The first column `odroid-0:/sharedDir` specifies the NFS server node and the NFS directory it exports, the second column specifies the local NFS directory, and the third column specifies the file system type.

Once the NFS server and the NFS clients have been configured, the students check that the NFS works properly. To do so we ask the students to access all the compute nodes, write a file in the shared directory from that node, and then check that all the files are visible by all the nodes.

F. Software Environment for Parallel Programming

Once all the hardware and the system software is set up, the last step to have a fully operational Odroid cluster is to install the software environment for the two standard parallel programming modes used in HPC: MPI and OpenMP. In order to understand the behaviour and the performance of the parallel programs we also install Extrae [15], an instrumentation library to transparently trace MPI and OpenMP programs, and Paraver [16], a trace visualiser and analyser that will allow students to understand the execution of parallel applications. Both are freely available at the Barcelona Supercomputing Center tools website [17].

OpenMP is available by default with the GNU compiler (`gcc` in the case of C), so the students do not need to install it manually. In contrast, the MPI packages are not installed by default, so the students must do it. Among the multiple implementations available for MPI we opt to use MPICH, which only requires to install the `mpich` package in every compute node.

To install Extrae the students first have to install all the required packages in all the compute nodes: `libtool`, `automake`, `m4`, `perl`, `libunwind-dev`, `libxml2-dev`, `binutils-dev`, and `libiberty-dev`. Then the students manually configure, compile and install the Extrae library in the NFS shared directory. To do this they use the commands

```
./bootstrap; ./configure --with-mpi=/usr
--with-unwind=/usr --without-dyninst
--without-papi --disable-parallel-merge
--prefix=/sharedDir/extrae; make -j 8; make
install. To use Extrae we provide the students with
scripts that automatically enable the tracing of their parallel
programs and post-process the generated trace so it is ready
to be analysed with Paraver.
```

Finally, the students download Paraver directly as a pre-compiled binary for `x86_64` architectures, which can be executed straight away in the head node.

REFERENCES

- [1] "<https://wiki.odroid.com/odroid-xu4/odroid-xu4>."
- [2] "MPI: A Message-passing Interface Standard. Version 3.1. June 2015."
- [3] "OpenMP Application Program Interface. Version 4.5. November 2015."
- [4] E. Ayguade and D. Jimenez-Gonzalez, "An approach to task-based parallel programming for undergraduate students," *Journal on Parallel and Distributed Computing*, vol. 118, no. P1, pp. 140–156, 2018.
- [5] E. Ayguade, L. Alvarez, and F. Banchelli, "OpenMP: what's inside the black box?" *Peachy Parallel Assignments (EduHPC 2018)*, 2018.
- [6] S. Williams, A. Waterman, and D. Patterson, "Roofline: An insightful visual performance model for multicore architectures," *Commun. ACM*, vol. 52, no. 4, pp. 65–76, Apr. 2009.
- [7] "<https://www.samsung.com/semiconductor/minisite/exynos/products/mobileprocessor/exynos-5-octa-5420>."
- [8] P. Greenhalgh, "big.LITTLE Processing with ARM Cortex-A15 & Cortex-A7," *ARM White paper*, 2011.
- [9] S. Holt, A. Meaux, J. Roth, and D. Toth, "Using inexpensive microclusters and accessible materials for cost-effective parallel and distributed computing education," *Journal of Computational Science Education*, vol. 8, no. 3, pp. 2–10, Dec. 2017.
- [10] D. Toth, "A portable cluster for each student," in *2014 IEEE International Parallel Distributed Processing Symposium Workshops*, 2014, pp. 1130–1134.
- [11] J. C. Adams, J. Caswell, S. J. Matthews, C. Peck, E. Shoop, D. Toth, and J. Wolfer, "The micro-cluster showcase: 7 inexpensive beowulf clusters for teaching pdc," in *Proceedings of the 47th ACM Technical Symposium on Computing Science Education*, 2016, pp. 82–83.
- [12] S. J. Matthews, "Teaching with parallella: A first look in an undergraduate parallel computing course," *Journal of Computer Sciences in Colleges*, vol. 31, no. 3, pp. 18–27, Jan. 2016.
- [13] "XSEDE. <https://www.xsede.org/>."
- [14] "iBlue Waters. <http://www.ncsa.illinois.edu/enabling/bluewaters>."
- [15] "Barcelona Supercomputing Center. Extrae User Guide Manual. Version 2.2.0. 2011."
- [16] V. Pillet, J. Labarta, T. Cortes, and S. Girona, "Paraver: A Tool to Visualize and Analyze Parallel Code," in *Proceedings of WoTUG-18: transputer and occam developments*, vol. 44, no. 1, 1995, pp. 17–31.
- [17] "<https://tools.bsc.es/downloads>."