

Mapping and Characterization of Applications in Heterogeneous Distributed Systems

Jaime Yeckle and Wilson Rivera
Computing and Information Sciences and Engineering
Parallel and Distributed Computing Laboratory
University of Puerto Rico, Mayaguez Campus
Mayagüez, PR 00681-9042
[jyeckle, wrivera}@ece.uprm.edu](mailto:{jyeckle, wrivera}@ece.uprm.edu)

Abstract

The problem of mapping tasks and communications onto multiple machines and networks in a heterogeneous computing environment has been shown to be NP complete. Therefore, the development of heuristic techniques to find near-optimal solutions is required. Many different types of mapping heuristics have been developed in recent years. However, selecting the best heuristic to use in any given scenario remains a difficult problem. Moreover, it is not possible to make a general mapping in a heterogeneous computing environment. In this paper we propose to characterize classes of applications with the objective of predicting their behavior. Using the insight provided by the characterization, we achieve a more realistic mapping for specific applications.

1 Introduction

The steady decrease in cost and increase in performance of commodity workstations and personal computer have made it increasingly attractive to use clusters of such systems as compute servers instead of high-end parallel supercomputers [16]. Due to the rapid advance in performance of commodity computers, when such clusters are upgraded by addition of nodes, they become heterogeneous. The issue of effective mapping of applications onto such heterogeneous clustered systems is therefore of great interest. Several research studies have addressed this problem [5,42,8,32,41,1,11,10,37]. However the problem is NP complete [9,22], therefore, the development of heuristics techniques is required. Many factors make it difficult to select the best technique of mapping. These factors can be described as follows:

1. A distributed computing environment has conflicting requirements [14]: (i) While minimizing interprocessor communication tends to assign the entire computation to a single processor, load balancing tries to distribute the computation evenly among the processors. (ii) While real-time constraints require many processors as possible to maximize parallel execution, the precedence relationships limit parallel execution. (iii) The saturation effect suggests the use of fewer processors since inefficiency increases with the number of processors.
2. When one heuristic technique is presented and evaluated in the literature, typically, different assumptions are made about the underlying target platforms making comparisons problematic. Similarly, different assumptions about application models complicate comparisons. In addition, the algorithms should take into account characteristics of processors, network architecture and applications. Regarding applications, for example, the existing algorithms only consider the size property.

In this paper, we propose to characterize classes of applications with the objective of predicting their behavior. In order to obtain this characterization of applications, we consider observations and statistical methods. Using the insight provided by the characterization, we propose a more efficient and realistic mapping for specific applications.

The rest of the paper is organized as follows. Section 2 begins with the description of different mapping schemes and taxonomy for describing matching and scheduling heuristics for heterogeneous computing systems. Section 3 presents the methodology for characterizing applications and the use for heterogeneous distributed systems. The paper

concludes with comments about the proposed strategy and its extensions in section 4.

2 Taxonomy of Mapping

Mapping includes assigning (matching) each task to a machine and ordering (scheduling) the execution of the tasks on each machine [4]. The mapping problem is extremely difficult to solve and generally intractable [3,26]. Even the simplified subproblems constructed from the original mapping problem by imposing a variety of constraints still fall in the class of NP hard problems. The difficulty of solution varies with the inclusion or exclusion of preemption, the number of parallel processors, precedence constraints, etc.

We now classify the various strategies for multiprocessor scheduling, task mapping, and resource allocation under a common, uniform set of terminology [7]. The Figure 1 shows the structure of the hierarchical portion of the taxonomy. The strategies can be classified as being either static or dynamic.

Static Mapping versus Dynamic Mapping: In the static mapping case the entire information regarding the processes in the host system, as well as the processes involved in a job, is assumed to be available a priori [41,28,37,32]. On the other hand, in the dynamic mapping a more realistic assumption is used, that is very little a priori knowledge is available about the resource needs of a process. In the static case, a decision is made for a process image before it is ever executed, while in the dynamic case no decision is made until a process starts.

Optimal versus Suboptimal: In the case that all information regarding the state of the system as well as the resource needs of a process are known, an optimal assignment can be made based on some criterion function [19,29,35,30]. Examples of optimization measures are minimizing total process completion time and maximizing utilization of resources in the systems.

In the case that these problems are computationally infeasible, suboptimal solutions may be obtained [28,33,40]. Within the realm of suboptimal solutions to the mapping problem, the heuristic algorithms represent the category of static algorithms that make a realistic assumption about a priori knowledge concerning process and host system characteristics. The most distinguishing feature of heuristic schedulers is that they make use of special parameters, which affect the system in indirect ways.

Optimal and Suboptimal Approximate Techniques: Regardless a static solution is optimal or

suboptimal approximate, there are four basic categories of task allocation algorithms, which can be used to arrive at an assignment of processes to processors:

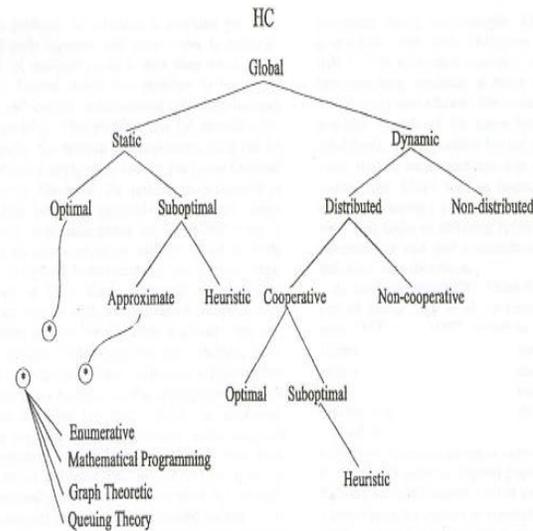


Figure 1. Taxonomy of mapping

- Solution space enumeration and search
- Graph theory
- Mathematical programming
- Queuing theory

Distributed Versus Nondistributed: Distributed mapping means that the work involved in making decisions should be physically distributed among the processors [15]. On other hand, nondistributed means whether the responsibility for the task of global dynamic scheduling physically resides in a single processor [31].

Cooperative versus Noncooperative: Cooperative means that all mechanisms which involve cooperation between distributed components and Noncooperative whether the individual processors make decisions independent of the actions of the other processors

Many algorithms have been published addressing the problem of matching and scheduling, where several simplifying assumptions are common. Orduña [32], for example, describes a mapping scheme assuming all the network switches are attached to the same number of workstations, the workstations are uniprocessors, and only one process is mapped to each processor. These assumptions nevertheless determine system performance.

Another simplifying assumption is made in [13]. In this paper mapping is modeled with forward flow only. Programs with dynamic structures are not considered. Also, Ahmad and Kwok [1] compared several algorithms for scheduling task graphs. The algorithms have different assumptions: bounded and unbounded number of processors and clusters, task duplication based scheduling and arbitrary processor network scheduling.

In the previous discussion, most of the approaches focus primarily on specific mapping strategies for particular multiprocessor architectures. Some approaches intend to take advantage of hardware characteristics, such as the interconnection network of architectures. Others take advantage of characteristics of processors and very few papers take advantage of characteristics of applications. Also exist other approaches, such as [10] which intent to generalize mapping and scheduling. By large, these approaches assume homogeneous conditions for many characteristics of applications or architecture of hardware. Therefore, a general mapping to all applications and characteristics of hardware and network is difficult to obtain. We propose a specific mapping considering hardware issues in the context of specific applications. To achieve this, we propose a characterization of the applications with the objective to predict behavior.

3 Characterization of Applications

In this section, we provide a new approach to mapping in heterogeneous distributed system considering the characterization of applications.

3.1 Characteristics of Applications

The characteristics of the applications (which can be tasks or subtasks) can be defined as follows:

Application size: This mean if a task contains subtasks or not.

Application type: Types of applications to be mapped.

Communication patterns: This mean to the sources and destination subtasks for each data item to be transferred.

Data availability: The time at which input data needed by a subtask or output data generated by a subtask can be utilized varies in relation to subtask start and finish times: (a) is data available (to be forwarded) before a subtask completes, and (b) can a subtask begin execution before receiving all of its input data? As an example, a clustering non-uniform assumes that a subtask cannot send data to other

waiting subtasks until it completely finishes executing.

Deadlines: the applications have deadlines.

Execution time model: Most mapping techniques require an estimate of the execution time of each application on each machine. The two choices most commonly used are probabilistic and deterministic modeling. Probabilistic modeling uses a probability distribution for application execution times when make mapping decisions [2, 27]. Deterministic modeling uses a fixed (or expected) value [17], e.g., the average of ten previous executions of an application.

Task heterogeneity: For each machine exist different tasks with different properties (e.g., probability distribution) which make the execution times different.

Multiple versions: the applications have multiple versions that could be executed For example, an application that requires an FFT might be able to perform the FFT with either of two different procedures that have different precisions, execution times, and resource requirements.

Priorities: Priorities are generally assigned by the user (within some allowed range), but the relative weight given to each priority are usually determined by another party (e.g., a system administrator). Priorities and their relative weightings are required if the mapping strategy is preemptive.

Task profile: Task profiling specifies the types of computations occurring in an application based on the code for the task (or subtask) and the data to be processed [20, 34]. This information may be used by the mapping heuristic, in conjunction with analytical benchmarking to estimate task (or subtask) execution time.

Temporal distribution: It is mean static applications (the complete set of tasks to be mapped known a priori), dynamic applications (the tasks arrive in a real-time, non-deterministic manner), or it can be a combination of the two

3.2 Characterization

The predictions of behavior in applications are based on past observations. These observations can be obtained using simulations or in real time. The survey in real time has not advantage because real applications of interest might run for long periods of time and it is not feasible to perform a statistically significant number of experiments. In addition, using real resources makes it difficult to explore a wide variety of resource configurations. Finally, variations in resource load over time make it difficult to obtain repeatable results. Simulation is then the most viable

approach to effectively evaluate scheduling algorithms.

Currently a few simulation packages are available [18, 38, 12, 39], but they are not targeted to the simulation of distributed applications for the purpose of evaluating scheduling algorithms. These tools are often very complete and sophisticated but too complex and low-level for our purpose. Simgrid [6] provides core functionalities that can be used to build simulators for the study of application scheduling in distributed environments. In this survey we consider Simgrid to perform simulations. In order to make simulations, we use a mapping algorithm which consider heterogeneity of processors and network.

3.3 Statistical method

In order to obtain the characterizations of applications using past observations we propose use statistical methods. This approach offers a number of advantages. First, a statistical method can compensate for many different factors, without requiring a distinct model for each of the different machine architectures. Second, statistical estimates will improve with time, as the number of previous observations increases. Finally, statistical schemes can be made to be computationally efficient, making them practical for use at run time. One potential criticism of statistical schemes is the need of a large number of past observations to obtain accurate estimates.

One important method to consider is nonparametric regression. Nonparametric regression has the advantage of being able to estimate equation of model, as a function of several parameters, without any knowledge of the function itself. Since we make no assumptions on the functional form, this prediction scheme does not require any knowledge of either the task of the target architecture, making it applicable in a very general sense.

Most of the previous work for heterogeneous distributed computing which mixing past observations and statistical methods are generally for execution time estimation. The SmartNET heterogeneous scheduling tool offers statistical execution time estimation technique, but no details of its implementation have yet been published [25]. Iverson [23, 24] present nonparametric regression technique, for execution time estimation in heterogeneous distributing computing. Other authors [21, 36] use techniques based on Bayesian decision theory, but these techniques are difficult to implement.

For prediction behavior of applications defined in this paper, we consider a function $p(\delta)$ where δ is a vector of parameters of the application such that size,

etc. While the estimation algorithm does not know any details about the functional form of $p(\delta)$, it does have a set of n previous observations of applications $\{(y(\lambda)_i, \delta_i)\}_{i=1}^n$, where $y(\lambda)_i$ is the prediction behavior of application for the vector δ and λ is a function that depends of time, bandwidth, etc. These observations are assumed to contain some amount of random error ε , such that

$$y(\lambda)_i = p(\delta_i) + \varepsilon_i$$

Thus, the goal of the problem is, given the function $p(\delta)$, to obtain an estimate of the characteristics of this application, using the set of previous observations.

3.4 Integration

The scheme will operate in the following manner. A set of previous observations is maintained by the algorithm. Using this set of observations, the prediction behavior of application in each potential machine can be estimated, and then matching and scheduling algorithm can use these estimates to make a realistic and efficient mapping decision. After the task of the application is complete, the new parameters are added to the set of observations, to be used to improve future predictions. By storing past observations, the estimation algorithm is able to improve its estimates over time.

4 Conclusions and future work

In this paper, we have presented a review of different techniques used to perform matching and scheduling in heterogeneous distributed systems. We propose a characterization of specific applications with the objective to obtain a forecast behavior. The predictions of behavior in applications are made base on a combination of observations, simulation and modeling through statistical methods. The result can be utilized to generate new mapping techniques more efficient and realistic. Our future work will be concerned with the development of mapping algorithms based on the new approach. We envision a set of mapping strategies and tools to reach peak performance of applications on heterogeneous distributed systems.

Acknowledgement

This work was supported by the UPRM-NSF PRECISE Project (NSF-EIA 99-77071) and the

References

- [1] I. Ahmad, K. Kwok and M. Wu. Analysis, evaluation, and comparison of algorithms for scheduling task graphs on parallel processors. *IEEE Transactions on Parallel and Distributed Systems*, 1996.
- [2] R. Armstrong, D. Hensgen, and T. Kidd, The relative performance of various mapping algorithms is independent of sizable variances in run-time predictions," 7th Heterogeneous Computing Workshop (HCW '98), Mar. 1998.
- [3] D. Bernstein, M. Rodeh, and I. Gertner. On the complexity of scheduling problems for parallel/pipelined machines. *IEEE Transactions on computers*, C-38(9):1308-1313, September 1989.
- [4] T. D. Braun, H. J. Siegel, N. Beck, L. Boloni, M. Maheswaran, A. I. Reuther, J. P. Robertson, M.D. Theys, and B. Yao, "A Taxonomy for Describing Matching and Scheduling Heuristics for Mixed-Machine Heterogeneous Computing Systems," 17th IEEE Symposium on Reliable Distributed Systems, Oct. 1998, pp. 330-335.
- [5] T. Braun, H. Siegel and A. Maciejewski. Static mapping heuristic for task with dependencies, priorities, deadlines and multiple versions in heterogeneous environments. *IEEE proceedings of the international parallel and distributed processing symposium* 1530-2075. Feb. 2002.
- [6] H. Casanova. Simgrid: a toolkit for the simulation of application scheduling. Computer Science Department University of California, San Diego, Sept. 2001.
- [7] T. L. Casavant and J. G. Kuhl. A taxonomy of scheduling in general purpose distributed computing systems. *IEEE Transactions on Software Engineering*, 14(2):42-45, February 1988.
- [8] B. Cirou, and E. Jeannot. Triplet: a Clustering Scheduling Algorithm for Heterogeneous Systems. IEEE Symposium on Reliable Distributed Systems, Oct. 2001.
- [9] E. G. Coffman, Jr., ed., *Computer and Job-Shop Scheduling Theory*, John Wiley & Sons, New York, NY, 1976.
- [10] V. Chaudhary, and J. K. Aggarwal. A generalized scheme for mapping parallel algorithms. *IEEE Transactions on Parallel and Distributed Systems*, 1998.
- [11] S. Chen, L. Xiao and X. Zhang. Adaptive and virtual reconfigurations for effective dynamic job scheduling in cluster systems. *Proceedings of the 22nd International Conference on Dist. Comp. Systems*. March 2002.
- [12] J. Davis, M. Goel, C. Hylands, B. Kienhuis, E. Lee, J. Liu, X. Liu, L. Muliadi, S. Neuendorffer, J. Reekie, N. Smyth, J. Tsay, and Y. Xiong. Overview of the Ptolemy Project. Technical report ERL Technical report UCB/ERL No M99/37, Dep. EECS, University of California, Berkeley, July 1999.
- [13] V. Dixit-Radiya and D. Panda. Clustering and Intra-Processor scheduling for explicitly parallel programs on distributed-memory systems. *IEEE Transactions on Parallel and Distributed Systems*, August 1994.
- [14] K. Efe. Heuristic models of task assignment scheduling in distributed systems. *Computer*, 15:50-56, June 1988.
- [15] P. H. Enslow Jr. What is a "distributed" data processing system? *Computer*, 11(1):13-21, January 1978.
- [16] Ian Foster and Carl Kesselman. *The Grid: BluePrint for a New Computing Infrastructure*. Morgan Kaufmann Publishers, 1998
- [17] R. F. Freund, M. Gherrity, S. Ambrosius, M. Campbell, M. Halderman, D. Hensgen, E. Keith, T. Kidd, M. Kussow, J. D. Lima, F. Mirabile, L. Moore, B. Rust, and H. J. Siegel, Scheduling resources in multiuser, heterogeneous, computing environments with SmartNet," 7th Heterogeneous Computing Workshop (HCW '98), Mar. 1998.
- [18] P. Fishwick. *Simulation Model Design and Execution: Building Digital Worlds*. Prentice Hall, 1994.
- [19] A. Gabrielian and D. B. Tyler. Optimal object allocation in distributed computer systems. In *Proc. Int. Con. on Distributed Computer Systems*, pages 84-95, May 1984.
- [20] A. Ghafoor and J. Yang, "Distributed heterogeneous supercomputing management system," *IEEE Computer*, Vol. 26, No. 6, Jun. 1993.
- [21] C. Hou and K. Shin. Load sharing with consideration of future task arrivals in heterogeneous distributed real time systems. *IEEE Transactions on Computer*, 43(9):1076-90, Sept. 1994.
- [22] O. H Ibarra and C. E. Kim, "Heuristic Algorithms for Scheduling Independent Tasks on Nonidentical Processors," *Journal of the ACM*, vol. 24, no 2, Apr. 1977, pp. 280-289.
- [23] M. Iverson and G. Follen. Run statistical estimation of Task execution times for

- heterogeneous distributing computing. *Proceedings of HPDC*, May 1996.
- [24] M. Iverson. Statistical prediction of Task execution times through analytic benchmarking for scheduling in a heterogeneous environment. *IEEE Transactions on Computer*, vol 48 Nro 12, December 1999.
- [25] T. Kidd, D. Hensgen, L. Moore, R. Freund, D. Charley, M. Halderman, and M. Janakiraman, "Studies in the Useful Predictability of Programs in a Distributed and Homogeneous Environment," The Smartnet Home Page, <http://papaya.nosc.mil:80/SmartNet/>, 1995.
- [26] E.L. Lawler, J.K. Lenstra, and A.H.G.R. Kan. Recent developments in deterministic sequencing and scheduling: a survey. In M. A. H. Dempster et al., editor, *Deterministic and Stochastic Scheduling*. D. Reidel publishing company, 1982.
- [27] Y. A. Li and J. K. Antonio, "Estimating the execution time distribution for a task graph in a heterogeneous computing system," 6th Heterogeneous Computing Workshop (HCW '97), Apr. 1997.
- [28] V. M. Lo. Heuristic algorithms for task assignment in distributed systems. In *Proc. Int. Conf. on Distributed Comp. Systems*, C-37(11):1384-1397, November 1988.
- [29] P. Y. R. Ma, E. Y. S. Lee, and J. Tsuchiya. A task allocation model for distributed computing systems. *IEEE Transactions on Computers*, C-31(1):41-47, January 1982.
- [30] L. M. Ni and K. Hwang. Optimal load balancing for a multiple processor system. In *Proc. Int. Con. on Parallel Processing*, pages 352-357, 1990.
- [31] J. Ousterhout, D. Scelza, and P. Sindhu. Medusa: An experiment in distributed operating system structure. *Communications ACM*, 23(2):92-105, February 1980.
- [32] J. Orduña and F. Duato. A new task mapping technique for communication-aware scheduling strategies. *IEEE Transactions on Parallel and Distributed Systems*, Set. 2001.
- [33] C. C. Price and S. Krishnaprasad. Software allocation models for distributed computing systems. In *Proc. Int. Conf. on Distributed Comp. Systems*, pages 40-48, May 1984.
- [34] H. J. Siegel, M. Maheswaran, and T. D. Braun, "Heterogeneous distributed computing," *Encyclopedia of Electrical and Electronics Engineering*, J. Webster, ed., John Wiley & Sons, New York, NY, to appear.
- [35] C. C. Shen and W. H. Tsai. A graph matching approach to optimal task assignment in distributed computing systems using a minimax criterion. *IEEE Transactions on Computers*, C-34(3):197-203, March 1985.
- [36] K. Shin and C. Hou. Design and evaluation of effective load sharing in distributed real time systems. *IEEE Transactions Parallel and Distributed Systems*, 5(7):704-19, July. 1994.
- [37] K. Taura and A. Chien. A heuristic algorithm for mapping communicating tasks on heterogeneous resources *IEEE Transactions on Parallel and Distributed Systems*, Set. 2000.
- [38] <http://www.threadtec.com/>. (Viewed September 2002).
- [39] A. Terzis, K. Nikoloudakis, L. Wang, and L. Zhang. IRL-Sim: A general-purpose packet level network simulator. In *Proceedings of the 33rd Annual Simulation Symposium*, Apr. 200. To appear.
- [40] A. M. VanTilborg and L. D. Wittie, "Wave scheduling – Decentralized scheduling of task forces in multicomputers," *IEEE Trans. Comp.*, vol. C-33, no 9, pp 835-844, Sept. 1984.
- [41] V. Yarmolenko, J. Duato, D. k. Panda, and P. Sadayappan. Characterization and Enhancement of Static Mapping Heuristic for Heterogeneous Systems. Technical report OSU-CIISRC-02/00-TR07, Dept. of computer science, Ohio State University, Feb. 2000.
- [42] A. Zomaya and T. Yee-Hwei. Observations on Using Genetic Algorithms for Dynamic Load Balancing. *IEEE Transactions on Parallel and Distributed Systems*, Set. 2001.