

Estimating Execution Time of Distributed Applications

Maciej Drozdowski*

Institute of Computing Science, Poznań University of Technology,
ul. Piotrowo 3a, 60-965 Poznań, Poland
Maciej.Drozdowski@cs.put.poznan.pl

Abstract. In this work we consider the problem of estimating execution time of distributed applications. The main difficulty stems from the communication delays and shared nature of the computing and communication media. A simple method taking into account communications and concurrency of computations is proposed to estimate the execution time of a distributed application. The proposed technique is evaluated in a sequence of experiments. It is demonstrated that this method is feasible. Other interesting consequences are also outlined.

Keywords: performance evaluation, communication delays, clusters of workstations, scheduling.

1 Introduction

Measurement of the execution time of a computer application is fundamental to assess its efficiency. This observation applies also to parallel programs. Contemporary parallel applications are executed in shared and spatially distributed environments. Consequently, parallel programs are exposed to communication delays and coexistence with many applications competing for the resources of the computer system. Hence, a method of measuring execution time of a distributed application in a shared environment would be very desired.

The need for efficient allocation of the activities (programs) to the resources (processors, communication channels) gave rise to many scheduling models in parallel processing (cf. monographs: [1,3,5,6]). These deterministic scheduling models assume initial knowledge of the parallel application parameters such as processing times or communication delays. In this work we consider a method of obtaining these parameters, the precision and soundness of the final model.

The study presented here has a practical origin. It has been observed that while reporting the efficiency of parallel applications run in shared environments, the influence of the competing programs is difficult to be dealt with. The indications of the computation efficiency tend to be unreliable. Thus, a practical method of eliminating such influence is needed.

* This research was partially supported by KBN grant.

The above considerations boil down to a question: How to measure distributed application execution time? In parallel applications run in a shared and distributed system, communication delays include influence of other communications. Likewise, computation time may include influence of other user applications.

Consider several ad hoc answers to the above question. Firstly, we can use *astronomical time* (wall time). This seems to be an ideal solution, but only in a completely dedicated system. Not only other users are not allowed to use our computers, our network segment must be completely isolated from the outside world, but also the software suite running on our computers must be fully controllable. These requirements are realistic in dedicated systems. However, it is far more difficult to fulfil them in geographically scattered and multiplatform systems. As a second solution, it is possible to measure processor times and system times consumed on all computers by our program. Still, this method does not take into account communication delays and precedence constraints between the operations executed on various processors. We conclude that a "good" method of measuring real execution time of parallel applications should:

1. account for the order of the program operations,
2. respect communication delays,
3. eliminate influence of other users and applications sharing the computing environment,
4. be simple conceptually and in implementation.

Satisfying the above requirements may be hard due to distributed nature of the application. We discuss it in the last section. However, we believe that a satisfactory estimation can be obtained by the method presented in this paper.

2 Proposition of a Solution

Without loss of generality let us assume that communications involve two different computers connected by some kind of network. Thus, we exclude for the time being broadcast and internal communications. We will use term *processor time* to denote both the processor time consumed by the program executed in the user space and by the system on behalf of the application.

In the scheduling theory it is common to represent activities, programs or projects by a precedence constraints graph. Let $G(V, A)$ denote a directed acyclic graph with set V of events representing fork or join of the program control, and set A of arcs representing communication or computation operations (cf. Fig.1). This method of representing precedences between the operations is called *activity on arc* in the deterministic scheduling theory [1]. Graph $G(V, A)$ represents control flow in the distributed application. Let every arc has a weight equal to the duration of the operation it represents. Hence, arc j representing computations has weight a_j equal to the processing (i.e. computation) time, while arc k corresponding to the communication has weight c_k equal to the communication time. The length of the longest path in graph G represents the execution

off-line as astronomical time. If $f(v)$ is measured in an unloaded or dedicated network, then the communication time estimated from $f(v)$ for known v , will retain character of a dedicated network even if the real communication took place in shared and loaded environment. In this way influence of the external communications can be eliminated. Thus, labeling communication arcs in graph G with the amount of data transferred is equivalent to labeling it with the communication time. The procedure of estimating the execution time is applied after the application termination on the basis of the recorded logs.

3 Experimental Evaluation

The method presented in the previous section has been tested in a series of experiments conducted on various applications and environments. The results are preliminary because the experiments were not always performed in the same way. We believe, however, that the main conclusions are valid despite that.

The goals of the evaluation were to verify stability of the reported running time in changing environment conditions, and to investigate coincidence of the results with the running time in dedicated environment. Satisfying the former condition is necessary for drawing conclusions on the behavior and the performance of distributed applications in shared and changing environment. The latter condition is imposed by the need for eliminating influence of the competing applications.

The first experiments (experiments 1) were performed in a cluster of four heterogeneous Sun workstations. The communication medium was a single segment 10Mb Ethernet. The second group of experiments (experiments 2) were performed in a cluster of five homogeneous PCs: Pentium 166MMX, 32 MB of RAM, Windows NT 4.0, connected by a single segment 10Mb Ethernet. The third series experiments (experiments 3) were done on six PCs: Pentium 200MMX, 64MB of RAM memory, with Linux operating system (Red Hat 6.0, kernel 2.2.5), interconnected by 100Mb single segment Ethernet. PVM message passing library was used in all experiments. The applications were: distributed search for a pattern in a text file (experiments 1), distributed multiplication of two matrices (experiments 2), distributed computation of Julia set image (experiments 3).

In the search for a pattern application [8] a big data file was divided into chunks. The chunks of text were sent to the processors to seek for a text pattern. A ready processor applied for data to the master processor. The order of activating the processors and the assignment of the data pieces to the processors was not known a priori. Therefore, the structure of the control flow graph G was known only after the execution of the program.

Matrix multiplication application [7] consisted in computing a product of two square matrices. The logical processor interconnection was a mesh topology. The multiplied matrices were divided into stripes and sent along the columns of the mesh. The results were sent along the rows of the mesh. In this application the control flow graph is known a priori. Nondeterminism was possible only in the order of final receiving the results by the master processor.

The third application [2] was a distributed computation of Julia set image. The image of the Julia set (a.k.a. Mandelbrot fractal) is showing the convergence of a sequence $z_0 = 0, z_{n+1} = z_n^2 + c$ of complex numbers. The convergence of z_n was mapped from a square $[-2,2] \times [-2,2]$ of parameter c values to a 1000×1000 bitmap. Slave processors analyzed convergence of z_n in a submesh determined by the coordinates of the opposite corners (lower-left and upper-right). Work was distributed to the processors in equal chunks representing some number of full horizontal lines in the bitmap. After initial distribution of the work the first processor requesting work also received new load first from the work distributor (master). Hence, the structure of the control flow graph was not known a priori as the order of sending work and receiving results from the slaves was not fixed.

To restore the control flow graph for all operations the consumed processor time, receiver/sender identifiers, amount of transferred data were logged. In experiments 1 communication time for the transferred amount of data was calculated using Lagrangean interpolation. In the remaining experiments communication time was approximated by a linear function.

The first test of the method consisted in verifying the coincidence of the astronomical execution time in a dedicated system, with the execution time calculated by our procedure in the system under varying load. These times should be equal, or at least, the difference should be constant. Example results of this test in experiments 1 are shown in Fig.2. Horizontal axis in Fig.2 represents system load reported by the operating system, along the vertical axis relative difference between the astronomical and the estimated times is shown. As it can be seen the difference is less than 18% and it is distributed in 6% range. Thus, the difference is small and stable.

In the second test stability of the method was examined. The changes of the execution time of the same application estimated by our method were recorded under various computer loads. The reported value should be constant under all loads. In experiments 3 the changes of the load were caused artificially by running a simple program with one infinite loop containing one integer add operation. The results of this test are presented in Fig.3 for experiments 3. On the vertical axis time is reported, on the horizontal axis the load is shown in the units of the number of running loading programs. The upper line is the astronomical execution time. Obviously, this time increases with the load. The lower line is the execution time estimated by our method. It is constant as expected.

Interesting results were obtained in experiments 2 (cf. Fig.4). The difference between the estimated execution time and the astronomical time was stable in the central range of the multiplied matrices sizes. However, the difference was non-negligible. It was observed, that the astronomical time comprises execution time of the operating system services and runtime environment (e.g. `pvmd`). These overheads are not included in the consumed processor time of our application. We concluded that astronomical time was longer than it should be. During the execution of our application, the operating system and other services consumed, e.g., 20% of the processor time. Thus, the astronomical time inevitably includes contribution of the operating system overheads. In order to have a fair compar-

ison the astronomical time should be decreased by approx. 20%. Alternatively, the computational arcs in control flow graph G should be lengthened approx. $1/(1 - 0.2)$ times. The results of such a correction in the estimated time are shown in Fig.4.

Influence of the communication time to processing time ratio has been also tested (cf. Fig.5). It was observed that with growing contribution of the processing time in the overall execution time, the difference between the estimated and astronomical times diminishes. This can be caused by decreasing contribution of the communication delays. It can be concluded that estimation of communication time is an important source of the error in our method.

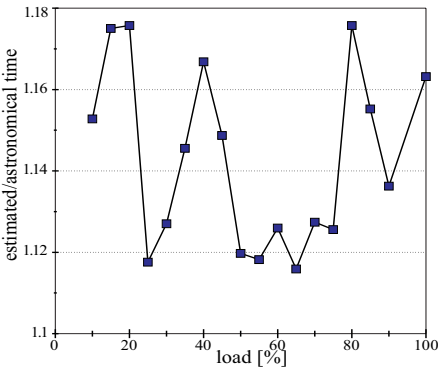


Fig. 2. Differences in the estimated execution time under varying load. Astronomical time measured in unloaded system. Experiments I.

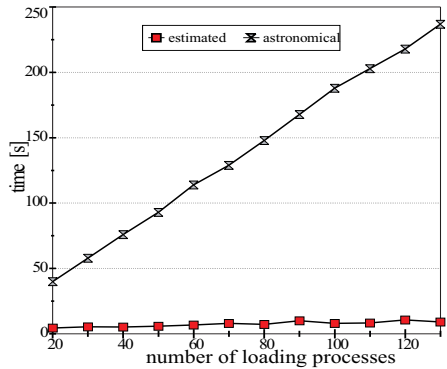


Fig. 3. Stability of the estimated execution time under varying load. Astronomical time changes with increasing load. Experiments III.

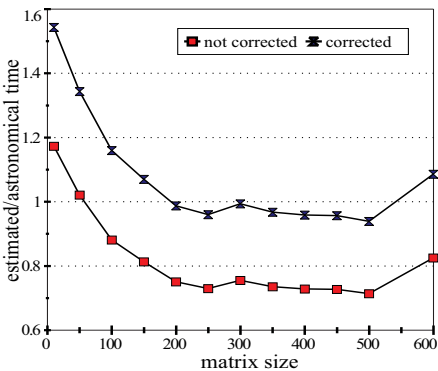


Fig. 4. Influence of the standard system load. Experiments II.

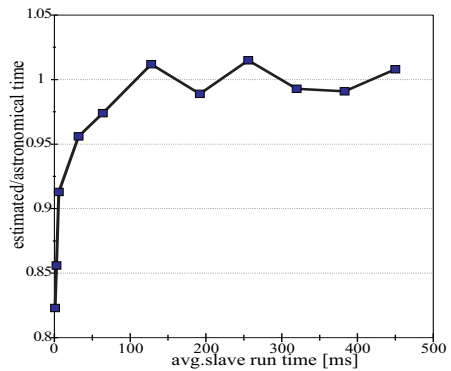


Fig. 5. Influence of the processing time duration. Experiments I.

4 Discussion

The of the experiments show that the method is rational. The calculated estimates are stable and close to the execution time in dedicated environment.

There are, however, limitations in drawing conclusions on the basis of the calculated estimates. In many cases the structure of control flow graph is non-deterministic, and can change, e.g., in different load conditions. Therefore, it may be unjustified to extend the results from the previous runs to the future ones.

There are potential sources of error in our method. The logging procedure may change the load of the computers. Consequently, the durations of the program operations and the structure of the control flow graph may be different than if the execution logs were not recorded. The results of our experiments show, however, that on average this influence is not big. Our method uses procedures of the operating system to measure consumed processor time. It is an open question how reliable these services are. Errors can be introduced also in the estimation of the communication delay. We observed that with the increasing speed of the communication networks and the processors the contribution of the linear dependence of communication time on the volume of transferred data diminishes. The nonlinear phenomena arising from communication initiation, buffering and reception become dominating. Note, that in all our tests these operations were implemented in software.

The results of experiments 2 (cf. Fig.4) lead us to yet another confusing observation that our method is able to eliminate influence of other application completely. This means that also operating system and runtime environment can be excluded from the estimation. Still, it would be difficult to use a computer system without such a software layer.

The technique we presented can be extended in various directions. It would not be difficult to include interprocess communications done on the same processor. Furthermore, broadcast messages can be logged and included in the control flow graph. In our considerations we assumed full interconnection (clique) between the processors, which is hardly ever the case. As a result the messages of the same application may overlap in time and compete with each other for the communication medium. This leads to the extension of the communication operations. Yet, such an extension can be calculated if the communication network topology and the medium access protocol are known. For example in Ethernet, when two (or more) messages share the communication channel, the bandwidth decreases proportionately to the number of competing messages. Practical verification of these extensions can be a subject of future research.

5 Conclusions

A new method of estimating execution time of a distributed application in a shared environment has been presented. The results of the experiments demonstrate method usability. This technique can be applied in performance profilers. Still, many fundamental questions remain open.

References

1. J. Błażewicz, K. Ecker, E. Pesch, G. Schmidt, J. Węglarz, *Scheduling Computer and Manufacturing Processes*, Springer, Berlin, 1996.
2. B. Burba, B. Figas, L. Wasyluk, Pomiar czasu wykonania aplikacji w środowisku rozproszonym, B.Sc. thesis, Institute of Computing Science, Poznań University of Technology, 1999.
3. P. Chretienne and C. Picoleau. Scheduling with communication delays: A survey. In P. Chretienne, E.G. Coffman Jr., J.K. Lenstra, and Z. Liu, editors, *Scheduling theory and its applications*. J. Wiley, 1995.
4. T.H. Cormen, C.E. Leiserson, R.L. Rivest, *Introduction to Algorithms*, Massachusetts Institute of Technology, 1990.
5. E.G. Coffman Jr. (editor). *Computer and job-shop scheduling theory*. Wiley & Sons, New York, 1976.
6. M. Drozdowski, *Selected problems of scheduling tasks in multiprocessor computer systems*, Series: Monographs, No.321, Poznań University of Technology Press, Poznań, (1997), (see also <http://www.cs.put.poznan.pl/~maciejd/h.ps>).
7. R. Janasiak, Pomiar rzeczywistego czasu wykonania aplikacji w środowisku rozproszonym, M.Sc. thesis, Institute of Computing Science, Poznań University of Technology, 1998.
8. R. Maciejewski, Szacowanie czasu wykonania aplikacji w środowisku rozproszonym, M.Sc. thesis, Institute of Computing Science, Poznań University of Technology, 1997.