

A Task Allocation Tool for Multicomputers

L. Hluchý, M. Dobrucký, D. Dobrovodský

Institute of Computer Systems, Slovak Academy of Sciences
Dúbravská cesta 9, 842 37 Bratislava, SLOVAKIA
e-mail: upsyhluc@savba.sk

Abstract: This paper describes the static and dynamic task allocation tool under PVM environment for distributed memory parallel systems. For the static allocation the objective function is used to evaluate the optimality of the allocation of a task graph onto a processor graph. Together with our optimization method also augmented simulated annealing and heuristic move exchange methods are implemented. For dynamic task allocation the semidistributed approach was designed based on the division of processor network topology into independent and symmetric spheres. The task allocation tool is controlled by user window interface - allocation toolbox.

Keywords: distributed memory parallel systems, static and dynamic allocation, multicomputer, load balancing.

1 Introduction

Optimal planning of parallel program execution in distributed environment of a multicomputer on the basis of message passing solves the response speed. The theory of optimal allocation comes out from the assumption that the program execution time depends upon uniform load of the processors and upon interprocessor communication minimization. In this paper, our attention is concentrated on the diffusion method for static allocation .

For static allocation a parallel program can be represented by a task graph $S(V,E)$, vertices $V=\{1,2,\dots,q\}$ represent the program tasks and edges E represent mutual communication. The topology of available technical resources can be represented by processor interconnection graph $H(C,E_p)$, vertices of which, $C=\{1,2,\dots,p\}$, represent processors, and edges E_p represent communication channels or links. We assume the homogenous structure of the multicomputer system. The communication cost depends on message size and on distance between source and destination. Interprocessor distance d_{ij} between processors i and j is given by routing rules, i.e. as minimal number of hops from processor i into the processor j . Clearly, if $d_{ij}=0$ then $i=j$. Hardware graph detection includes methods for multiprocessor system inspecting and diagnostics (seeking connections and node parameters in the structure). Many optimization methods have been created to solve the static allocation problem, [1 - 9]. Dynamic task allocation for parallel and distributed computing is an important goal to be achieved for various applications. The purpose of dynamic task allocation is to increase the system throughput in a dynamic environment, which can be done by balancing the utilization of computing resources and minimizing communication overheads among processors during run time. Dynamic task allocation approaches have been studied by numerous researcher [11 - 13]. This paper has 7 sections. Sections 2 and 3 describe the specification of cost function for static allocation and the diffusion method for optimization of cost function. In section 4 is described dynamic task allocation. In section 5 is described the implementation of static task allocation under PVM. In section 6 is

Funded by Commission of the European Communities DG12 under contract number CIPA-CT93-0251 and Grant Agency for Science of Slovak Academy of Sciences

described the implementation of dynamic task allocation under PVM. In section 7 is described the task allocation toolbox.

2 The Cost (Objective) Function Definition

To specify an appropriate optimization goal it is necessary to create a cost function, which provides a realistic evaluation of the communication and computation overhead. For the given graphs S and H , allocation A , the general form of the cost function CF can be built as a combination of two parts - so called **vertex cost function** F_{vertex} , **edge cost function** F_{edge} :

$$CF(A, t) = F_{vertex}(A, t) + F_{edge}(A, t)$$

Variable t is duration of the iteration step of the allocation algorithm. $F_{vertex}(M, t)$ expresses the effect of computation loads so that individual physical processors would be loaded in a uniform way. $F_{edge}(M, t)$ expresses the effect of communication volume on allocation so that allocation should ensure the minimum of external (interprocessor) communications.

3 The Diffusion Method for Static Allocation

The principle of our allocation approach can be represented by the following procedure [10]. In the starting condition all tasks are located at the root node of the multicomputer network. The tasks are transferred by centrifugal force $f1$ (following from the requirement of uniform processor load), against which centripetal force $f2$ is acting, trying to keep the processes with mutual communication as close to each other as possible. This way, uniform load is obtained and this method can be considered as an improvement of the "pure" load balancing method. The resulting sum of forces is a vector with components for each task stored in a node and directed to each communication link. A task and a direction with maximum (positive) value is chosen. The following relation is the definition of changes in the cost function, where $\alpha(t)=const.$ and function $g(t)$ hides the weight coefficient:

$$\Delta CF(M) = \alpha(t) \cdot f1(t) + g(t) \cdot f2(t)$$

where t is the duration of iteration step.

The total number of iterations must be sufficient for location stabilization. This number is obtained from the experiments and is a trade-off between the time consumption and the result accuracy. We estimate this number from q - the number of tasks to be mapped. Standard- nominal value was estimated as $\frac{7q}{2}$.

The function $g(t)$ that we use was obtained from programming experiments in the transputer environment for some applications. Our experiments aimed at looking for suitable face of function $g(t)$ (linear, nonlinear, increasing and decreasing intervals,...) and the number of iterations. The function $g(t)$ is suggested to have the following properties: During the first n iterations it can be kept near zero or increases slowly from zero. The reason is that all tasks are located at the same processor at the start of allocation, and there is a need to keep the communication force low to enable task distribution. In the n -th iteration step $g(t)$ reaches the value where the influence of $f1$ and $f2$ is comparable. Then $g(t)$ decreases to zero, and thus the influence of communications becomes weaker. The algorithm can be implemented in the distributed environment as well.

4 The Dynamic Task Allocation

In the recent years the research of dynamic task scheduling in the parallel computer systems brought numerous different approaches. They vary from centralized to fully distributed ones. We have studied some of them from both categories. Centralized methods with a process responsible for making scheduling decisions (scheduler) running on the central node can result in exploitation of its whole computing capacity. If the set of nodes is too large and load information fluctuate rapidly, the communication overhead can overload the part of the interconnection network surrounding the central node. In the distributed approaches it is difficult to keep the node load information topical in the time. Messages with these information can increase the overall communication overhead. To avoid collisions when different schedulers are trying to allocate a task on the same node is time consuming.

We have chosen the compromise between centralized and distributed approach as is proposed in the paper [12] where the semi-distributed method of load balancing is described. The set of processors is divided into independent symmetric regions, called spheres. In each sphere, a processor equidistant from all other processors is selected as the scheduler for that sphere. Cumulative load information for each sphere is exchanged among the independent schedulers. Using knowledge about load of local nodes and information from other schedulers, the scheduler can place a requested task on the appropriate local node or transfer it into a less loaded sphere. Dividing the set of nodes into spheres considers granularity of the tasks. In the problems with small granules, the average lifetime of the process is short. Tasks arise and extinct more frequently than those in the problem with big granules. So the sphere served by one scheduler should be smaller in a way to avoid overload of the central node.

5 The Implementation of Static Allocation in PVM

The allocation tool we have developed up to now is implemented under the UNIX operating system as distributed parallel program running under PVM (Parallel Virtual Machine ver.3.3). It can be run in the background, because all the inputs are supplied via command-line parameters. It uses the SW-graph file as the input for the description of the parallel application to be mapped and creates an output file in which the result (the location vector) is stored. Our tool supports also various other allocation methods (simulated annealing [1], Boillat [4], heuristic move exchange [2] ... some of them have sequential version only).

6 The Implementation of Dynamic Task Allocation in PVM

The Parallel Virtual Machine (PVM) environment has a built-in support for dynamic task allocation. Using library function `pvm_regrm()` a task can register itself as a Resource Manager (RM). After registration every request for spawning a new PVM task is sent to the RM instead of the PVM daemon to which RM is registered. RM can decide where to spawn a new task with regard to the load balancing.

The application responsible for load balancing presented by us consist of four types of processes (Fig. 1). The main process has the user interface and is designed for a determination of spheres, a creation of schedulers and a control of whole system for load balancing. The Scheduler is a process designed to make task placement decisions. The Agent

running on the each node collects information about average load. Based on these information Scheduler is updating its database. RMs running on each node catch spawn requests from tasks and invoke the scheduler to decide where the tasks should be placed. Requests to create a new task are sent to the RMs running on destination nodes. They are responsible for the task creation on their nodes. To achieve global load balanced, schedulers inform each other about average load of the sphere.

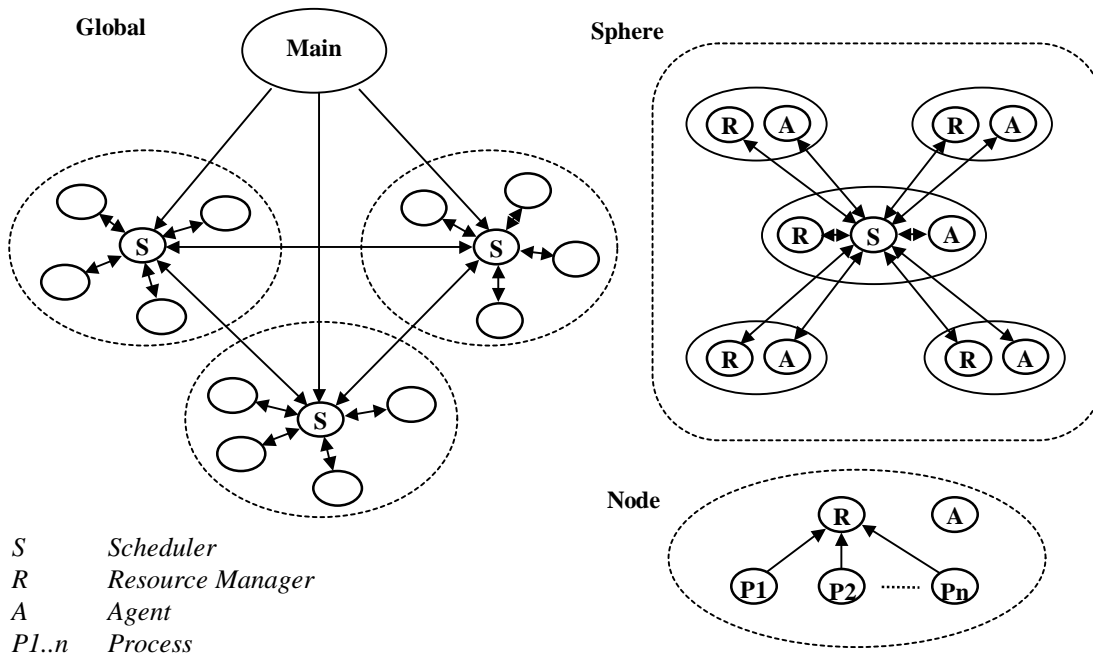


Fig. 1: Block scheme of dynamic task allocation.

7 The Task Allocation Toolbox

We have implemented the task allocation toolbox in the window Motif-like environment. In Figs 2, 3 examples of the toolbox window are given.

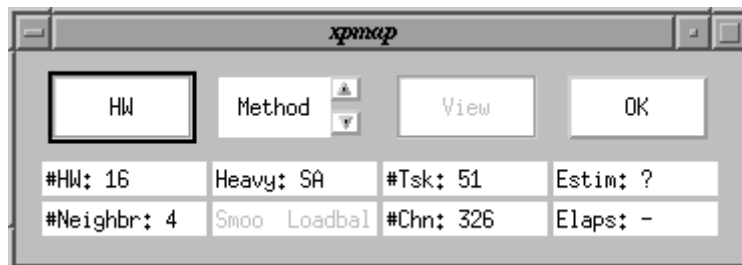


Fig. 2: Starting point - the user program (with 51 parallel tasks and 326 channels among them) is to be allocated onto parallel hardware machine with 16 processors (hypercube)

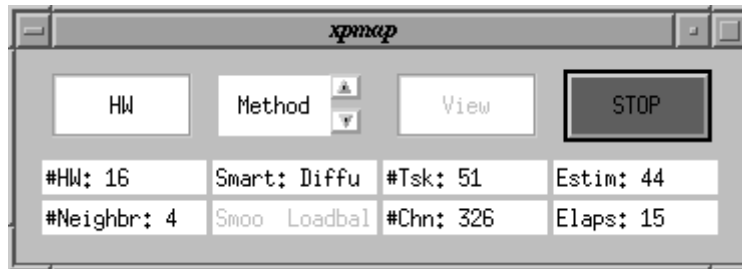


Fig. 3: Running phase - the Diffusion method was chosen.
This is the 15th iteration from 44.

When the task allocation toolbox (TATB) is started it has to read info about SW-graph (description of the application parallel program - task graph) and info about parallel hardware (hardware topology). Then the TATB window appears.

In the bottom part of the TATB window the status parameters are given :

- #HW - the number of PVM-hosts (size of the HW-graph)
- #Neighbr - the maximum of the processor neighbors - diameter of the HW-graph
- "Heavy: SA" - example of the chosen method
- "LoadBal" - indicator of 2-state switch for the dynamic load-balancing (dynamic task allocation)
- #Tsk - the number of parallel tasks to be mapped (SW- graph size)
- #Chn - the number of communication channels between parallel tasks
- "Estim" - estimated [time] consumption (it depends on the chosen method and SW-HW-graphs sizes)
- "Elaps" - progress in the [time] spending on the allocation

The execution button "**OK**" is available only if the allocation process is allowed to be started. Its face can be changed either to recommend further activities (that should be done, i.e. "HW?" if no HW- graph info was found in input file) or "Stop", "Wait" or "Restart". The allocation can be interrupted and started again, for example with the less time-consuming allocation method.

View: Our allocation tool can show the result progress (i.e. imbalance percentage, communication costs, time consumption, number of iterations ...) and also animates it.

Method: The user can change the allocation method. Allocation methods are grouped into 3 groups according their quality and also time consumption in the **quick** (fast simple greedy methods suitable for large tasks), **smart** (our diffusion method, Boillat method and method derived from our diffusion method suitable for a medium size tasks) and **heavy** (simulated annealing and "heuristic move-exchange", more time consuming but with the best achieved quality suitable for smaller tasks).

HW: After starting, our tool reads info about SW-graph and HW-graph. Activating this button shows the "pvm:" - number of the live PVM - nodes.

8 Conclusion

In this paper we proposed the cost function and diffusion method for the static task allocation for distributed memory parallel computers. Implementation is done in the distributed form under PVM. Together with the diffusion method also augmented simulated annealing, “heuristic move-exchange” method and several methods for faster allocation (greedy methods) were implemented. In our task allocation tool we have implemented dynamic task allocation based on the semi-distributed approach. We suppose two phase optimization for allocation of task graph into distributed memory parallel computers: in the first phase static task allocation is done (before compilation) and then during run time dynamic task allocation is executed which balances the utilization of computing resources and minimizes communication overheads among processors. For user we provide the task allocation toolbox in which he has possibilities to choose a static allocation method and switch on the dynamic task allocation. Implementation is performed on the generic UNIX operating system (Solaris 2.4, IRIX 5, Linux 1.2.x) and PVM 3.3.

In the future we want to verify our task allocation tool on the basis of concrete applications. So far randomly generated graphs were used for testing the task allocation tool.

References

- [1] KVASNIČKA, V. - POSPÍČHAL, J. - BISKUPIČ, S.: Task Allocation Problem Solved by Augmented Simulated Annealing. *Central European Journal for Operation Research and Economics*, in print.
- [2] SELVAKUMAR, S. - MURTHY, S.R.C.: An efficient heuristic algorithm for allocation parallel programs onto multicomputers. *Microprocessing and Microprogramming*, Vol.36, 1992/93, pp. 83-92.
- [3] BOKHARI, S.H.: On the Allocation Problem. *IEEE Transactions on Computers*, Vol. C-30, 1981, No. 3, pp. 207-214.
- [4] BOILLAT, J.E. - ISELIN, N. - KROPP, P.G.: MARC: A Tool for Automatic Configuration of Parallel Programs. In: Welch, P. et. al. (Eds.): *Transputing 91*, IOS Press 1991, pp. 311-329.
- [5] van LAARHOVEN, P.M.J. - AARTS, E.H.L.: *Simulated Annealing: Theory and Applications*. Reidel, Dordrecht, 1987.
- [6] BOLLINGER, S.W. - MIDKIFF, S.F.: Heuristic Technique for Processor and Link Assignment in Multicomputers. *IEEE Transaction on Computers*, Vol. 40, 1991, No.3, pp. 325-333.
- [7] BIANCHINI, R.P. - SHEN, J.P.: Interprocessor Traffic Scheduling Algorithm for Multiple-Processor Networks. *IEEE Transactions on Computers*, Vol.36, 1987, No. 4, pp. 396-409.
- [8] LEE, S.Y. - AGARWAL, J.K.: A Allocation Strategy for Parallel Processing. *IEEE Transactions on Computers*, Vol.36, 1987, No. 4, pp. 433-442.
- [9] COROYER, C. - LIU, Z.: Effectiveness of Heuristics and Simulated Annealing for the Scheduling of Concurrent Tasks-An Empirical Comparison. In: Bode, A.-Reeve, M.-Wolf, G. (Eds.): *Proceedings PARLE'93*. Springer 1993, pp. 452-463.
- [10] HLUCHÝ, L. - DOBRUCKÝ, M. - DUDÁK, M.: Solving Method for Optimal Load Balancing and Communication Minimisation. In: Plander, I. (Ed.): *Proceedings of the Sixth International Conference on Artificial Intelligence and Information - Control Systems of Robots*, World Scientific 1994, pp. 297-302.
- [11] CHANG H.W. D., OLDHAM W.J.B.: Dynamic Task Allocation Models for Large Distributed Computing Systems. *IEEE Trans. on Parallel and Distributed Systems*, Vol. 6., No. 12, December 1995, pp. 1301-1315.
- [12] AHMAD I., GHAFOR A.: Semi-Distributed Load Balancing for Massively Parallel Multicomputer Systems. *IEEE Trans. on Software Engineering*, Vol. 17, No. 10, October 1991, pp. 987-1004
- [13] LIN H.Ch., Raghavendra C.S.: A Dynamic Load - Balancing Policy with a Central Job Dispatcher (LBC). *IEEE Trans. on Software Engineering*, Vol. 18, No. 2, February 1992., pp. 148-158.