



The Ministry of Education and Science of the Russian Federation

Lobachevsky State University of Nizhni Novgorod

Computing Mathematics and Cybernetics faculty

The competitiveness enhancement program
of the Lobachevsky State University of Nizhni Novgorod
among the world's research and education centers

Strategic initiative

“Achieving leading positions in the field of supercomputer technology
and high-performance computing”

Introduction to MPI

Lecture 11. Parallel Computation Modeling and Analysis

Nizhni Novgorod

2014

Lecture_11_. Parallel Computation Modeling and Analysis

The analysis of parallelism efficiency is a crucial point in the development of parallel algorithms for solving complicated research and engineering problems. Parallelism efficiency analysis is, as a rule, the evaluation of the computation process speedup (reducing the time needed for solving a problem). Forming the speedup estimation may be carried out for selected computational algorithm (the efficiency estimation of parallelizing a specific algorithm). Another important approach may be the construction of the maximum possible speedup estimation for the solution of a certain problem type (the efficiency estimation of the best parallel approach for solving a problem).

In this lecture we will describe the computation model as an “operations-operands” graph, which can be used for the description of the existing information dependencies in selected algorithms of problem solving. We will also give the maximum possible parallelism efficiency estimations, which may be obtained as a result of the analysis of the existing computation models. The practical uses of the theory described here are given in the third part of the teaching materials.

11.1. Computation Model as “Operations-Operands” Graph

The model “operations-operands” graph can be used for the description of the information dependencies in selected algorithms of solving problems (see, for example, Bertsekas and Tsitsiklis (1989)). To simplify the problem we will assume that in constructing a model the periods of execution of any computational operations will be the same and will be equal to 1 (in some units of measurement). Besides we will assume that the data transmission among computing processors is carried out instantaneously without any time consumption (which may be quite true, for instance, if there is a common shared memory in a parallel computing system). The analysis of the parallel algorithm communication complexity is carried out in the next chapter.

Let us depict the set of the operations, carried out in the computational problem solution algorithm to be studied, and the information dependencies, which exist among the operations as an *acyclic oriented graph*

$$G = (V, R),$$

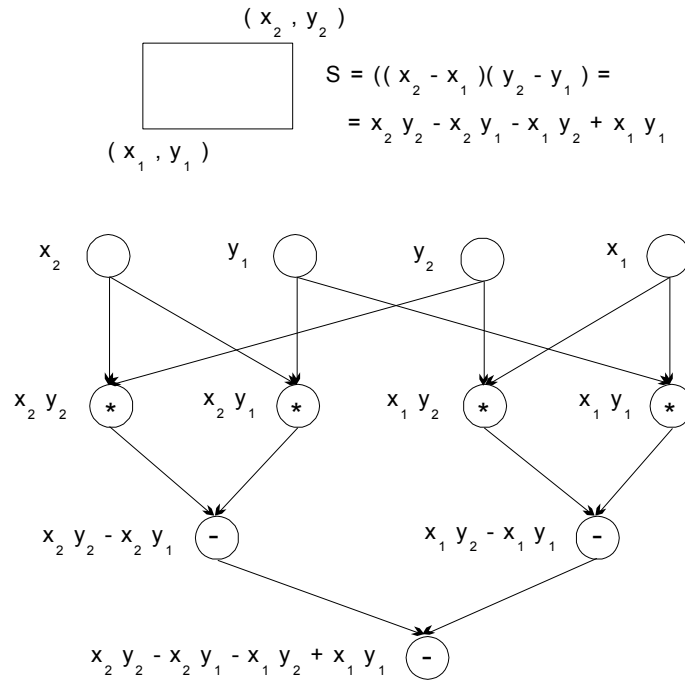


Figure 11.1. The Sample of computational model in the form of the “operations-
operands” graph

where $v = \{1, \dots, |V|\}$ is the set of graph vertices, which represent the algorithm operations being executed, and R is a set of graph arcs (in this case $r = (i, j)$ belongs to the graph only if the operation j makes use of the result obtained by execution of operation i). To illustrate this Figure 11.1 shows the graph of the algorithm used to calculate the area of the rectangle specified by the coordinates of its two opposite angles. As the given example shows, various computation schemes may be used and various corresponding computational models can be constructed to carry out the selected problem solution algorithm. As it will be shown later different computation schemes possess different capabilities of parallelizing. Thus the task of selecting the most suitable for parallel execution of a computational scheme algorithm can be set in constructing a computation model.

In the computational model of the algorithm under consideration the vertices without the incoming arcs may be used to assign the input operations, and the vertex without outgoing arcs may be used for output operations. Let us denote the set of graph vertices without input vertices as \bar{V} , and the diameter (length of maximum path) of the graph as $d(G)$.

11.2. The Scheme of Parallel Algorithm Execution

The algorithm operations, which do not have paths among them within the selected computation scheme, may be executed in parallel (for the computation scheme shown in Figure 11.1, for instance, first all the multiplication operations may be executed in parallel, and then the first two subtraction operations may be realized in parallel). A possible way to describe the parallel algorithm execution is given below (see, for instance, Bertsekas and Tsitsiklis (1989)).

Let p be the number of processors to execute an algorithm. Then to execute computations in parallel it is necessary to specify the set (*schedule*)

$$H_p = \{(i, P_i, t_i) : i \in V\},$$

where for each operation $i \in V$ the number of processor P_i used to execute the operation and the operation start time t_i are given. To make the schedule realizable it is necessary to meet the following requirements in specifying the set H_p :

- 1) $\forall i, j \in V : t_i = t_j \Rightarrow P_i \neq P_j$, i.e. the same processor must not be assigned to different operations simultaneously,
- 2) $\forall (i, j) \in R \Rightarrow t_j \geq t_i + 1$, i.e. all the necessary data must have been calculated before operation execution starts.

11.3. Evaluation of Parallel Algorithm Execution Time

The computation scheme of the algorithm G in combination with the schedule H_p may be considered as the model of the parallel algorithm $A_p(G, H_p)$, executed with the use of p processors. The time of parallel algorithm execution is determined by the maximum time value used in the schedule

$$T_p(G, H_p) = \max_{i \in V} (t_i + 1).$$

For the selected computation scheme it is desirable to use the schedule which provides the minimum algorithm execution time

$$T_p(G) = \min_{H_p} T_p(G, H_p).$$

The decrease of execution time may be provided by fitting the best computation scheme

$$T_p = \min_G T_p(G).$$

Estimates $T_p(G, H_p)$, $T_p(G)$ and T_p may be used as the time criteria in parallel algorithm execution. Besides to analyze the maximum possible parallelism it is possible to specify the estimate of the fastest algorithm execution

$$T_{\infty} = \min_{p \geq 1} T_p .$$

Estimate T_{∞} may be considered as the minimum possible time of the parallel algorithm execution if an unlimited number of processors are used (the concept of the computer system with the infinite number of processors usually called a *paracomputer* is widely used in the theoretical analysis of parallel computations).

Estimate T_1 defines the algorithm execution time if one processor is used and thus represents the execution time of the sequential version of problem solution algorithm. Constructing such an estimate is an important task in analyzing parallel algorithms, as it is necessary to evaluate the effect of the parallelism use (of speedup while solving the problem). It is evident that

$$T_1(G) = |\overline{V}| ,$$

where $|\overline{V}|$, as it has already been defined, is the number of vertices of the computational scheme G without the input vertices. It is important to note that if in determining the estimate T_1 we are limited to the consideration of only one selected problem solution algorithm and use the value

$$T_1 = \min_G T_1(G) ,$$

then the speedup coefficients obtained in accordance with the given estimate will characterize the efficiency of parallelizing the selected algorithm. To evaluate the efficiency of the parallel solution of the computational problems under consideration the time of the sequential solution must be evaluated with regard to various sequential algorithms, that is to use the value

$$T_1^* = \min T_1 ,$$

where the operation of minimum is taken over the set of all the possible sequential algorithms for a given problem.

We will consider the theoretical statements, which characterizes the properties of parallel algorithm execution time estimates (see Bertsekas and Tsitsiklis (1989)).

Theorem 1. The maximum path length of the algorithm computation scheme determines the minimum possible time of parallel algorithm execution, i.e.

$$T_{\infty}(G) = d(G) .$$

Theorem 2. Let there be a path from each input vertex for a certain output vertex in the algorithm computation scheme. Besides let the input power of the scheme vertices (the number of

incoming arcs) not exceed 2. Then the minimum possible time of parallel algorithm execution is limited from below by the value.

$$T_{\infty}(G) = \log_2 n ,$$

where n is the number of input vertices in the algorithm scheme.

Theorem 3. If the number of the used processors decreases, the algorithm execution time increases in proportion to the decrease of the number of processors, i.e.

$$\forall q = cp, \quad 0 < c < 1 \Rightarrow T_p \leq cT_q .$$

Theorem 4. For any number of the processors used the following upper estimate for parallel algorithm execution time is true:

$$\forall p \Rightarrow T_p < T_{\infty} + T_1 / p .$$

Theorem 5. The algorithm execution time comparable with the minimum possible time T_{∞} can be achieved if the number of processors is in the order of $p \sim T_1 / T_{\infty}$, to be precise,

$$p \geq T_1 / T_{\infty} \Rightarrow T_p \leq 2T_{\infty} .$$

If there are fewer processors, the time of algorithm execution cannot exceed the best computation time with the given number of processors more than twice, i.e.

$$p < T_1 / T_{\infty} \Rightarrow \frac{T_1}{p} \leq T_p \leq 2 \frac{T_1}{p} .$$

These theorems allow to form the basis for the following recommendations concerning the rules of parallel algorithm creation:

- 1) The graph with the minimum possible diameter must be used while choosing the algorithm computation scheme (see Theorem 1);
- 2) The efficient number of processors for parallel execution is determined by the value $p \sim T_1 / T_{\infty}$ (see Theorem 5);
- 3) The parallel algorithm execution time is limited from above by the values given in Theorems 4 and 5.

In order to specify the recommendations on the creating the schedule of parallel algorithm execution we will consider the proof of theorem 4.

The proof of the theorem 4. Let H_{∞} be the schedule for achieving the minimum possible execution time T_{∞} . For each iteration τ , $0 \leq \tau \leq T_{\infty}$, of the H_{∞} schedule execution the number of

operations carried out during the iteration τ will be written as n_τ . The schedule of the algorithm execution with the use of p processors may be constructed in the following way. We will divide the algorithm execution into T_∞ steps; at each step τ all n_τ operations, which were carried out during the iteration τ of the H_∞ schedule, must be carried out. The execution of these operations must be accomplished not more than in $\lceil n_\tau / p \rceil$ iterations with the use of p processors. As a result, the time of algorithm T_p execution may be evaluated the following way:

$$T_p = \sum_{\tau=1}^{T_\infty} \left\lceil \frac{n_\tau}{p} \right\rceil < \sum_{\tau=1}^{T_\infty} \left(\frac{n_\tau}{p} + 1 \right) = \frac{T_1}{p} + T_\infty.$$

The proof of the theorem offers a practical method of constructing the parallel algorithm schedule. First the schedule with no regard for the limitations of the number of used processors may be created (a paracomputer schedule). Then according to the scheme of the theorem derivation, the schedule for a finite number of processors can be constructed.

11.4. Parallel Algorithm Efficiency Characteristics

Speedup. This is a speedup obtained if a parallel algorithm is used for p processors in comparison to the sequential computations. It is determined by the value

$$S_p(n) = T_1(n) / T_p(n),$$

i.e. as the ratio of the problem solution time on a scalar computer to the time of parallel algorithm execution (value n is used for parameterization of computation complexity of the problem being solved and can be understood as, for instance, the amount of input problem data).

Efficiency. The efficiency of the processor utilization by the parallel algorithm in solving a problem is determined by the formula

$$E_p(n) = T_1(n) / (pT_p(n)) = S_p(n) / p$$

(the efficiency value determines the mean fraction of algorithm execution time, during which the processors are actually used for solving the problem).

Selecting the necessary parallel method of problem solving, it is very useful to estimate the computation cost, which is defined as the product of the parallel problem execution time and the number of the processor being used.

$$C_p = pT_p.$$

In this connection it is possible to define the concept of the *cost-optimal* parallel algorithm, which is defined as the method, the cost of which is proportional to the time of the best sequential algorithm execution.

To illustrate the introduced concepts in the next chapter we will consider a case of solving the problem of calculation of the partial sum for sequence of numerical values. Besides, in part 3 of the teaching materials these characteristics are used to estimate the efficiency of the considered parallel algorithm for solving the typical problems of computational mathematics.

11.5. Estimation of Maximum Attainable Parallelism

1. Amdahl's law. Maximum speedup obtainment may be hindered by the presence of sequential calculations in the computations being carried out, as the former cannot be parallelized. Let f be the part of the sequential calculations in the applied data processing algorithm, then, in accordance with Amdahl's law, the computation process speedup, if p processors are used, is limited by the value

$$S_p \leq \frac{1}{f + (1 - f) / p} \leq S^* = \frac{1}{f}.$$

Thus, for instance, if there are only 10% sequential instructions in the executed computations, the impact of parallelism use cannot exceed the tenfold data processing speedup. For the problem under consideration the computation of the sum of values for the cascade scheme the part of the sequential computations is $f = \log_2 n / n$. As a result, the value of the possible speedup is limited by the estimate $S^* = n / \log_2 n$.

Amdahl's law characterizes one of the most serious problems in the area of parallel programming (there are practically no algorithms without a certain part of sequential instructions). However, the part of sequential actions characterizes very often the sequential feature of the applied algorithms and does not characterize the possibility of parallel problem solution. As a result, the part of sequential computations may be decreased considerably if we choose methods that more appropriate for parallelizing.

It should be also mentioned that Amdahl's law is considered under the assumption that the part of sequential computation f is a constant value and does not depend on the parameter n , which defines the computational problem complexity. However, for a great number of problems the part $f=f(n)$ is a descending function of n . In this case the speedup for a fixed number of processors may be increased at the expense of increasing the computational complexity of the problem to be solved. This remark may be formulated as the statement that the speedup $S_p = S_p(n)$ is

the ascending function of the parameter n (this statement is often referred to as the *Amdahl's effect*). Thus, for example, for a problem under consideration – the computation of the sum of values – when a fixed number of processors p are used, the summarized data set may be subdivided into blocks of n/p size. Partial sums may be computed in parallel for the blocks first. Then these sums may be summarized with the help of the cascade scheme. The duration of the sequential part of the executed operations (minimum possible parallel execution time) is in this case

$$T_p = (n / p) + \log_2 p ,$$

that leads to the estimation of the sequential computation part as the value

$$f = (1 / p) + \log_2 p / n .$$

This expression shows that the sequential computation fraction f decreases with the increase of n . And in the limiting case we will obtain the ideal estimate of the maximum possible speedup $S^*=p$.

11.6. Analysis of Parallel Computation Scalability

The aim of parallel computation application is in many cases not only to decrease the computation execution time, but also to provide the possibility of solving more complicated variants of problem (such statements of the problem which cannot be solved if only uniprocessor computing systems are used). The parallel algorithm capability to efficiently use processors when the computation complexity increases is an important characteristic of the executed calculations. In this connection, the parallel algorithm is referred to as a scalable algorithm if with the increase of the number of processors it provides the speedup increase maintaining constant level of efficiency in processor use. A possible method to characterize the scalability properties is described below.

Let us assess the *total overhead* expenses, which take place in parallel algorithm execution

$$T_0 = pT_p - T_1 .$$

The total overhead expenses arise, as it is necessary to organize the interaction of processors. It is also necessary to fulfill some additional actions, synchronization of parallel computation and etc.

Making use of the previously introduced notation we can get new expressions for the time of solving the parallel problem solution and the speedup corresponding to it:

$$T_p = \frac{T_1 + T_0}{p}, \quad S_p = \frac{T_1}{T_p} = \frac{pT_1}{T_1 + T_0}.$$

With the use of the obtained relation the efficiency of the processor use may be expressed as

$$E_p = \frac{S_p}{p} = \frac{T_1}{T_1 + T_0} = \frac{1}{1 + T_0 / T_1}.$$

The latter expression shows that if the problem complexity is fixed ($T_1 = \text{const}$), then the efficiency will decrease if the number of processors increases at the expense of the total overhead costs T_0 . If the number of processors is fixed, the efficiency of processor used may be improved by the increase of the complexity T_1 of the problem being solved (it is assumed that with the increase of the complexity parameter n the total overhead expenses T_0 increase more slowly than the amount of computations T_1). As a result, if the number of processors increases, the necessary level of efficiency may be provided in the majority of cases by means of the corresponding problem complexity increase. In this connection the proportion of the necessary rates of calculation complexity increase and the number of processors being used becomes an important feature of parallel computations.

Let $E = \text{const}$ be the desirable efficiency level of the executed computations. Using the equation for the efficiency we may obtain

$$\frac{T_0}{T_1} = \frac{1 - E}{E} \quad \text{or} \quad T_1 = KT_0, \quad K = E / (1 - E).$$

The dependency $n = F(p)$ between the problem complexity and the number of processors generated by the latter relation is referred to as *isoefficiency function* (see Kumar et al. (1994)).

To illustrate this we will show the derivation of the isoefficiency function for the problem of summarizing numeric values. In this case

$$T_0 = pT_p - T_1 = p((n/p) + \log_2 p) - n = p \log_2 p$$

and the isoefficiency function looks as

$$n = Kp \log_2 p.$$

As a result, for instance, to provide the efficiency level $E = 0.5$ (i.e. $K = 1$) when the number of processors is $p = 16$, the number of summarized values must not be smaller than $n = 64$. If the number of processors is increased from p to q ($q > p$) it is necessary to increase the number n of

the summarized values $(q \log_2 q)/(p \log_2 p)$ times to provide the proportional speedup increase $(S_q/S_p)=(q/p)$.

11.7. References

Additional information on parallel computation modeling and analysis may be found in, for instance, Bertsekas and Tsitsiklis (1989). Useful information is also contained in Kumar et al. (1994), Quinn (2004).

The consideration of the academic problem of the numeric value sequence summation was carried out in Bertsekas and Tsitsiklis (1989).

For the first time Amdahl's law was stated in Amdahl (1967). Gustafson-Barsis's law was published in Gustafson (1988). The concept of isoefficiency was proposed in Grama et al. (1993).

A systematic discussion (for the time when the book was published) of the parallel computation modeling and analysis issues is given in Zomaya (1996).

11.8. Discussions

1. How is the "operations-operands" model defined?
2. How is the schedule for the distribution of computations among processors defined?
3. How is the time of parallel algorithm execution defined?
4. What schedule is optimal?
5. How can the minimum possible time of problem solving be defined?
6. What is a paracomputer? What can this concept be useful for?
7. What estimates should be used as the characteristics of the sequential problem solving time?
8. How to define the minimum possible time of parallel problem solving according to "operands-operations" graph?
9. What dependences may be obtained for parallel problem solving time if the number of processor being used is increased or decreased?
10. What number of processors corresponds to the parallel algorithm execution time (periods) comparable in the order with the estimates of minimum possible time of problem solving?
11. How are the concepts "speedup" and "efficiency" defined?
12. Is it possible to attain superlinear speedup?
13. What is the contradictoriness of the speedup and efficiency characteristics?
14. How is the concept of computation cost defined?
15. What is the concept of the cost-optimal algorithm ?
16. What does the problem of parallelizing a sequential algorithm of the numeric values summation lie in?

17. What is the essence of the summation cascade scheme? What is the aim of considering the modified version of the scheme?
18. What is the difference between the speedup and efficiency characteristics for the discussed versions of the summation cascade scheme?
19. What is the parallel algorithm of all the partial sums computation of a numeric value sequence?
20. How is Amdahl's law formulated? Which aspect of parallel computation does it allow to take into account?
21. Which algorithm is scalable? Give examples of methods with different level of scalability.

11.9. Exercises

1. Develop a model and evaluate speedup and efficiency of the parallel computations:

- For the problem of the scalar product of two vectors

$$y = \sum_{i=1}^N a_i b_i ,$$

- For the problem of choosing the maximum and minimum values for the given set of numeric values

$$y_{\min} = \min_{i \leq N} a_i, \quad y_{\max} = \max_{i \leq N} a_i ,$$

- For the problem of finding the mean value for the given set of numeric values

$$y = \frac{1}{N} \sum_{i=1}^N a_i .$$

2. Evaluate according the Amdahl's law the maximum attainable speedup for the problems given in 11.1

References

Amdahl, G. (1967). Validity of the single processor approach to achieving large scale computing capabilities. In AFIPS Conference Proceedings, Vol. 30, pp. 483-485, Washington, D.C.: Thompson Books.

Bertsekas, D.P., Tsitsiklis, J.N. (1989). Parallel and distributed Computation. Numerical Methods. - Prentice Hall, Englewood Cliffs, New Jersey.

Grama, A.Y., Gupta, A. and Kumar, V. (1993). Isoefficiency: Measuring the scalability of parallel algorithms and architectures. IEEE Parallel and Distributed technology. 1 (3). pp. 12-21.

Gustavson, J.L. (1988) Reevaluating Amdahl's law. Communications of the ACM. 31 (5). pp.532-533.

Kumar V., Grama, A., Gupta, A., Karypis, G. (1994). Introduction to Parallel Computing. - The Benjamin/Cummings Publishing Company, Inc. (2nd edn., 2003)

Quinn, M. J. (2004). Parallel Programming in C with MPI and OpenMP. – New York, NY: McGraw-Hill.