

Poznań University of Technology

**Multi-Installment Divisible Loads Scheduling
in Systems with Limited Memory**

J.Berlińska, M.Drozdowski, M.Lawenda

Research Report RA-07/08
2008

Institute of Computing Science, Piotrowo 2, 60-965 Poznań, Poland

Multi-Installment Divisible Loads Scheduling in Systems with Limited Memory

J.Berlińska¹, M.Drozdowski^{2*}, M.Lawenda³

1) Faculty of Mathematics and Computer Science,
Adam Mickiewicz University, Umultowska 87, 61-614 Poznań

2) Institute of Computing Science,
Poznań University of Technology, Piotrowo 2, 60-965 Poznań, Poland

3) Poznań Supercomputing and Networking Center,
Noskowskiego 10, 61-704 Poznań, Poland

Abstract

In this paper we study divisible load scheduling in systems with limited memory. Divisible loads are parallel computations which can be divided into independent parts of arbitrary sizes and processed in parallel on remote computers. The problem consists in distributing the load taking into account communication time, computation time, and limited memory buffers, such that the whole processing lasts as short as possible. The distributed system is a heterogeneous single level tree (star). The amount of memory at the processors participating in the computation is too small to accommodate the whole load at any moment of time. Therefore, the load is distributed in many small installments. Memory reservations have block nature, by which we depart from earlier models simplifying the memory management. We formulate this problem as a mixed nonlinear programming problem, and then propose two algorithms to solve it. The branch-and-bound algorithm is nearly unusable due to its complexity. Then, a genetic algorithm is proposed with more tractable execution times. We extensively analyze impact of various system parameters on the quality of the solutions. From this we infer on the nature of the scheduling problem.

Keywords: Parallel processing, scheduling, divisible loads, memory limitations, multiple installments.

1 Introduction

In this paper we study scheduling divisible loads in systems with limited memory sizes. Divisible load (DL) model represents parallel computations which can be divided into parts of arbitrary sizes, and the parts can be processed independently in parallel. These two simple assumptions on the nature of the parallel application have deep implications. Namely, the grains of parallelism are negligibly small because the load part sizes may be arbitrary. Since the parts can be processed independently, there are no data dependencies or other kinds of precedence constraints in the computation. Parallel processing of big

*Corresponding author, Maciej.Drozdowski@cs.put.poznan.pl. Research partially supported by Polish Ministry of Science and Higher Education.

volumes of data conforms with DL model. The processed data is generally called as *load*. Divisible load model originated in the late 1980ties [1, 4]. In [1] DL model has been applied to represent distributed computations in a network of workstations. In publication [4] a chain of intelligent sensors was considered. In both cases the problem was how to partition the computations such that the whole processing time is as short as possible. On one hand distributing the computations reduces processing time by employing additional computers. On the other hand, distributing the computations takes time. Hence, the problem is what load quantities should be sent to which processors. The mathematical models proposed in the early publications were easily tractable and boiled down to systems of linear equations. Later on, DL model has been applied to analyze performance of various computer network topologies, systems with parameters varying in time, limited memory sizes, to schedule computations in minimum monetary costs, and other. Overall, the divisible load theory (DLT) delivered a generic and versatile method of analyzing a broad class of parallel computations. Surveys of DLT can be found, e.g., in [2, 5, 11].

Scheduling divisible loads in systems with limited memory sizes was first considered in [9] where a heuristic called Incremental Balancing Strategy was proposed. It was assumed that all the processors always take part in the computation, the sequence of sending the load pieces to the processors is given, and that the whole load fits in the memory buffers of the computers. Furthermore, a simple linear model of communication delay was assumed. A more general affine communication time model including communication startup times was analyzed in [6]. A linear programming formulation was given which delivers optimum load partitioning under affine communication time model, and for a given sequence of load distributing. The problem of constructing optimum set of processors participating in the computation was shown to be **NP**-hard in [8]. A branch-and-bound (BB) algorithm and a bunch of heuristics has been proposed and experimentally evaluated in [8] for the problem of single-installment divisible load processing with arbitrary communication sequence, and arbitrary set of participating processors. In this paper we assume that the whole load is too big to store it in the memories of the computers at the same moment. Therefore, it is distributed and processed in many small pieces each of which fits in the computer memory. This organization of computations is called a *multi-installment* or *multi-round* processing. Multi-installment processing of divisible loads in systems with limited memory has been analyzed in [7]. An affine communication delay function was assumed, the set of processors taking part in the computation and the sequence of communicating with them was arbitrary. A branch-and-bound, and genetic (GA) algorithms have been proposed to solve this problem. However, memory management has been simplified in [7] to make the mathematical model tractable (we discuss it in more detail in the next section).

In this paper we study scheduling divisible loads in a star system (also called a single level tree) with limited memory buffers. We assume that the communication delay is an affine function of the amount of processing load. The set of processors taking part in the computation, and the sequence of sending load

chunks to them can be arbitrary and must be selected by the scheduling algorithm. Moreover, we assume that the memory reservations, and releases have realistic block nature (it is detailed in the next section). We express our problem as a mathematical programming formulation, and propose two algorithms to solve it: an optimization branch-and-bound algorithm (BB) which guarantees optimality of the solution, and an approximate genetic search algorithm (GA). Though it can be proved that the former algorithm (BB) delivers optimum solutions, it is practically unusable due to its complexity. The latter algorithm (GA) delivers good quality solutions only on average, but it is more applicable considering the execution time. We propose the GA not only to define yet another metaheuristic solving some combinatorial optimization problem, but also to gain some insight into the features of near-optimum solutions. An extensive computational study was conducted to analyze practical features of the scheduling problem important for processing large divisible computations.

The rest of the paper is organized as follows. In Section 2 we formulate the problem formally. In Section 3 methods of solving the problem are presented and discussed. We report on the results of the computational experiments, and the insights into the nature of the problem in Section 4. The last section is dedicated to the conclusions.

2 Problem Formulation

In this work we assume that each processing element is equipped with a CPU, some memory, and hardware front-end for managing network communications (e.g. NIC and DMA). The CPU and network hardware can work in parallel such that simultaneous computation and communication is possible. The words processing element, processor, and computer will be used interchangeably. We assume star interconnection. In the center of the star a processor P_0 called originator is located. The originator is connected to a set $\{P_1, \dots, P_m\}$ of processors. Initially the originator has some volume V of load to be processed. The load is sent directly from the originator to the processors. The star topology may represent a cluster of workstations connected via a local area network, a set of CPUs in SMP system sharing a bus, a distributed computation with the originator as a master, and the computers as workers in a grid environment. We assume that only processors P_1, \dots, P_m perform computations, and the originator does no computing. Were it otherwise, the computing capability of the originator can be represented as an additional processor. For simplicity of mathematical models we assume that the time of returning the results to the originator is negligible. The processor and its communication link to the originator are characterized by the parameters: A_i - computing rate (inverse of speed e.g. in seconds per byte), B_i - size of available memory (expressed, e.g., in bytes), C_i - communication rate (inverse of bandwidth), S_i - communication startup time (e.g. in seconds). The process of load distribution consists in sending pieces of the load to the processors for remote computation. Words installment, chunk, message, piece of load, communication will be used interchangeably. The transmission time of

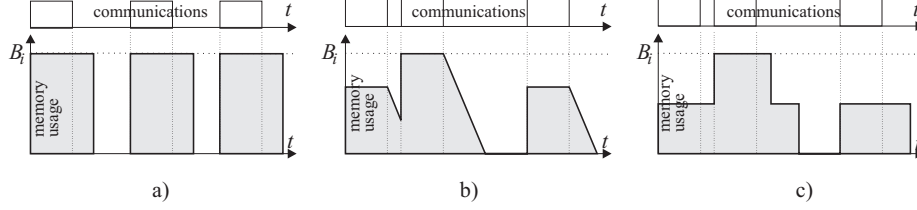


Figure 1: Memory management: a) each chunk uses whole buffer, b) memory gradually released, c) block memory releases.

a load chunk of size α (e.g. bytes) sent to processor P_i is $S_i + \alpha C_i$. The same amount of load is computed on P_i in time αA_i .

Now let us analyze memory management. We assume that memory is allocated from the operating system pool at the beginning of the communication comprising the load chunk, and it is released to the operating system after the end of computation on the load chunk. The size of the load which arrived at a processor may not exceed the amount of available memory. The simplest approach to modeling memory usage is to assume that only one chunk is held by a processor at a time (cf. Fig.1a). Hence, the chunk of size α_i may use all the available memory and a constraint $\alpha_i \leq B_i$ is imposed, for processor P_i . This approach was used in [8, 9]. Yet, it is insufficient in multi-installment processing when many messages may arrive at the processor: The load can be gradually uploaded, new and old buffers can be swapped without stopping the computations. Consequently, the load chunk sizes may interact with each other. In [7] it was assumed that the received load *together* shall not exceed memory size B_i . However, it was also assumed that memory is released to the operating system with very fine granularity equal to the load unit. Consequently, memory occupation was decreasing linearly during computations (cf. Fig.1b). With such a simplification load chunk sizes could be calculated using linear programming for a given communication sequence. Still, this way of releasing memory seems rather unusual. In this paper we assume that memory allocation and release have block nature (cf. Fig.1c). When a chunk of size α_j is about to arrive to a processor, a block of α_j load units is requested from the operating system. This block exists in the memory pool of the application until finishing computation on chunk j . On completion of chunk j a block of size α_j is released to the operating system. The sizes of coexisting memory blocks cannot exceed limit B_i . In other words, for each moment t , $\sum_{l \in \mathcal{H}(i,t)} \alpha_l \leq B_i$, where $\mathcal{H}(i,t)$ is the set of chunks received by P_i and not completed by time t . We will be saying that chunks simultaneously existing in memory buffer *overlap*. Thus, $\mathcal{H}(i,t)$ is a set of chunks overlapping at time t on P_i .

We will express our scheduling problem as a mixed nonlinear mathematic program. Let us introduce necessary assumptions and notation. The load is delivered to the processors in a sequence of communications. The sequence may be arbitrary, which means that some processors may be excluded from

Table 1: Notation used in formulation (1)-(15)

m	number of processors;
A_i, B_i, C_i, S_i	computing rate, memory size, communication rate, communication startup time of processor P_i , respectively;
σ	communication sequence (constant in (1)-(12));
$\sigma(i)$	index of the processor receiving the i -th chunk in (1)-(12);
n	length of communication sequence σ (constant in (1)-(12));
n_i	number of chunks sent to processor i (constant in (1)-(12));
$\rho(i, k)$	global number of the k -th chunk received by processor i , i.e. mapping from local chunk numbers on processor i to the global numbers (constant in (1)-(12));
M	a constant greater than schedule length, for example, $M \geq V(\max_{j=1}^m C_j + \max_{j=1}^m A_j) + n \max_{j=1}^m S_j$ (constant in (1)-(12));
α_i	the i -th chunk size (variable in (1)-(12));
T_{max}	schedule length (variable in (1)-(12));
t_i	time when sending of message i (global number) starts (variable in (1)-(12));
f_{ik}	time when processing of message k (local number) on processor i finishes (variable in (1)-(12)).
x_{ijk}	binary variable denoting if chunks j, k overlap on P_i (variable in (1)-(12)).
δ_{ij}	advancement of overlap interval introduced with chunk j on P_i (variable, used in alternative overlap encoding(13)-(15)).
z_{ij}	the last chunk overlapping with chunk j on P_i (derived from δ_{ij} , used in alternative overlap encoding(13)-(15)).
$n_{MIN} = \frac{V}{\max_i \{B_i\}}$	minimum number of load chunks
$m' \leq m$	the number of different used processors

the computations, while some other processors may receive the load more often than others. If the message is received by a processor without any load in the buffer, then computations start immediately after the end of communication. If the buffer already stores some unprocessed chunks, then the processor switches from computing one load chunk to the next one without idle time in the computations. Idle times may arise between the communications when processor memory occupation is maximum, and no new load may be uploaded to any processor. We assume that the sequence of processing the chunks on a given processor is the same as the sequence in which they were received. Let us assume that the sequence $\sigma = (\sigma(1), \dots, \sigma(n))$ of the communications to the processors is given, where $\sigma(i)$ is the index of the processor receiving the i -th chunk. The numbers of the load chunks as they are sent off the originator will be called global numbers. For simplicity of notation also local numbering of the chunks received by a certain processor will be used. Function $\rho(i, j)$ is a mapping from processor P_i local chunk number j to the global numbering. In the following mathematical program we want to express the fact that chunks simultaneously existing in a processor buffer do not exceed memory size. To formulate such a constraint we have to know which load chunks overlap, but this depends on the communication sequence, and chunk sizes which are unknown. Consequently, the sets of overlapping chunks are to be determined. Let x_{ijk} be a binary variable equal 1 if the j -th chunk on processor P_i overlaps with chunk k , and equal 0 otherwise. In other words, $x_{ijk} = 1$ means that P_i started receiving chunk k before computing the j -th chunk was finished. In Table 1 we summarize the notation introduced so far and the notation used in the following mathematical program. Our problem can be formulated in the following way.

minimize T_{max}
subject to

$$t_1 = 0 \quad (1)$$

$$t_i \geq t_{i-1} + S_{\sigma(i-1)} + C_{\sigma(i-1)}\alpha_{i-1} \quad i = 2, \dots, n, \quad (2)$$

$$f_{ik} \geq t_{\rho(i,k)} + S_i + C_i\alpha_{\rho(i,k)} + A_i\alpha_{\rho(i,k)} \quad (3)$$

$$i = 1, \dots, m, \quad k = 1, \dots, n_i,$$

$$f_{ik} \geq f_{i,k-1} + A_i\alpha_{\rho(i,k)} \quad (4)$$

$$i = 1, \dots, m, \quad k = 2, \dots, n_i,$$

$$f_{ij} \geq t_{\rho(i,k)} - (1 - x_{ijk})M \quad i = 1, \dots, m, \quad (5)$$

$$j = 1, \dots, n_i - 1, \quad k = j + 1, \dots, n_i$$

$$f_{ij} \leq t_{\rho(i,k)} + x_{ijk}M \quad i = 1, \dots, m, \quad (6)$$

$$j = 1, \dots, n_i - 1, \quad k = j + 1, \dots, n_i$$

$$x_{ijk} \leq x_{ilk} \quad i = 1, \dots, m, \quad j = 1, \dots, n_i - 1, \quad (7)$$

$$k = j + 2, \dots, n_i, \quad l = j + 1, \dots, k - 1$$

$$x_{ijk} \geq x_{ijl} \quad i = 1, \dots, m, \quad j = 1, \dots, n_i - 1, \quad (8)$$

$$k = j + 1, \dots, n_i, \quad l = k + 1, \dots, n_i$$

$$\alpha_{\rho(i,j)} + \sum_{k=j+1}^{n_i} x_{ijk} \alpha_{\rho(i,k)} \leq B_i \quad i = 1, \dots, m, j = 1, \dots, n_i \quad (9)$$

$$V = \sum_{i=1}^n \alpha_i \quad (10)$$

$$T_{max} \geq f_{in_i} \quad i = 1, \dots, m \quad (11)$$

$$x_{ijk} \in \{0, 1\} \quad (12)$$

Variables α_i in the above formulation define load partitioning resulting in minimum schedule length for the communication sequence σ . Inequalities (1), (2) determine the moment when sending of the i -th chunk starts. Constraints (3),(4) determine the earliest time moment f_{ik} when computation on chunk k of P_i finishes. Inequalities (5), (6) guarantee that according to the value of x_{ijk} processing of chunk j is, or is not, finished before starting message k . Only one of inequalities (5), (6) is active for some i, j, k . If $x_{ijk} = 0$ then chunk j should be finished before starting the k -th communication to processor i . In this case inequality (6) is active (and (5) is inactive) ensuring the requirement. If $x_{ijk} = 1$, then chunk j is still unfinished when message k is initiated. In this case inequality (5) is active ((6) is inactive) ensuring the overlap of chunks j, k . Inequalities (7) guarantee that if chunk j is not processed when chunk k arrives, then the chunks between j , and k are also unprocessed. Inequalities (8) ensure that if chunk j is finished before arriving of some chunk k , then j can no longer become unprocessed again. By inequalities (9) memory limits are observed. No load remains unprocessed by (10). Schedule length is not shorter than the completion time on any processor by constraints (11). Formulation (1)-(12) is a mixed quadratic mathematical program because we have binary variables (x_{ijk}), continuous variables ($\alpha_i, f_{ik}, t_i, T_{max}$), and multiplication of variables in constraints (9). This implies that program (1)-(12) is hard to solve even though the activation sequence σ is given. This is in sharp contrast with the complexity of memory management models used in [7, 8] for which linear programs were needed to calculate load partition for a given communication sequence σ . Thus, representation of block memory management, and chunk overlap made the mathematical model much more involved. Note that (1)-(12) is very general and may cover various scenarios of optimum memory management. For example, it is capable of representing a number of independent buffers of equal or different sizes swapped on the processors. Let us observe that for a given x_{ijk} the above formulation becomes a linear program (LP). This is a foundation of the solution methods proposed in the next section.

3 Algorithms

In the previous section we established that for a given communication sequence σ , and given values of variables x_{ijk} encoding chunk overlap, chunk sizes α_i can be calculated by an LP. Hence, our problem can be solved by a tandem of

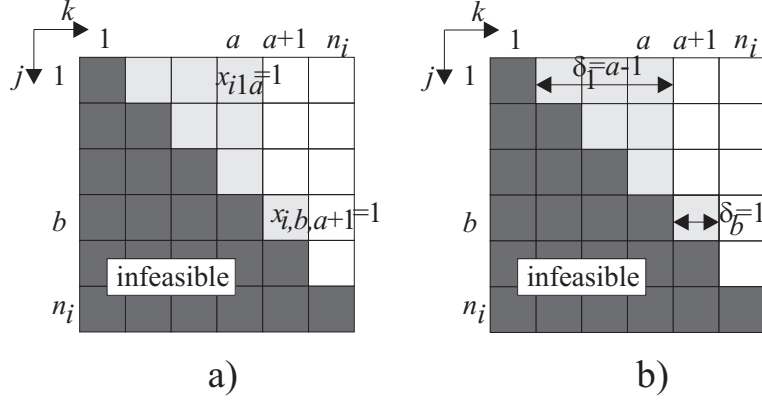


Figure 2: Overlapping encoding. a) by x_{ijk} , according to the (1)-(12), b) alternative encoding using overlap extensions. Black area - infeasible x_{ijk} , light gray - $x_{ijk} = 1$, white - $x_{ijk} = 0$.

methods. The first solves the combinatorial part of the problem by selecting activation sequence σ , and chunk overlap x_{ijk} . The second part solves the LP for the given σ and x_{ijk} . Below we present a convenient solution encoding method, and two methods working according to the above tandem rule.

3.1 A New Solution Encoding Method

In formulation (1)-(12) the overlapping of several chunks is determined by variables x_{ijk} . Since there are at most $n(n-1)/2$ such variables (when $\exists_i n_i = n$) the number of possible value assignments is at most $2^{n(n-1)/2}$. However, taking into account constraints (7),(8) a different way of encoding chunk overlapping is possible. For some chunk j sent to processor P_i it is possible to encode how many chunks received on P_i after j overlap with processing of chunk j . Furthermore, due to constraints (7) it is only necessary to encode by how many positions the front of overlapping is shifted ahead with each new message received by P_i (see Fig.2). Instead of $n_i(n_i-1)/2$ binary variables x_{ijk} we will use n_i integer variables δ_{ij} denoting by how many chunks the overlapping front is forwarded with chunk j on P_i . If the values of δ_{ij} , for $j = 1, \dots, n_i$, are given then constraints (5),(6) may be rewritten as follows:

$$f_{ij} \geq t_{\rho(i, z_{ij})} \quad \text{for } i = 1, \dots, m, j = 1, \dots, n_i - 1 \quad (13)$$

$$f_{ij} \leq t_{\rho(i, z_{ij}+1)} \quad \text{for } i = 1, \dots, m, j = 1, \dots, n_i - 1 \quad (14)$$

where $z_{ij} = \min\{n_i, \max\{z_{i, j-1} + \delta_{ij}, j + \delta_{ij}\}, \}$, and $z_{i1} = 1$ (cf. Fig.2). Here z_{ij} is the index of the last chunk which overlaps chunk j . If $z_{ij} = n_i$, only constraints (13) are necessary, and constraints (14) must be dropped because $\rho(i, j)$ is undefined for $j > n_i$. Furthermore, if $z_{ij} = n_i$ for some j then constraints (13) for $j' > j$ are not necessary, as they are enforced by (4). Analogously, if

$z_{ij} = j$, constraint (13) may be dropped as enforced by constraint (3). In the new encoding the bounds on the memory usage (9) are replaced with:

$$\sum_{k=j}^{z_{ij}} \alpha_{\rho(i,k)} \leq B_i \quad i = 1, \dots, m, j = 1, \dots, n_i \quad (15)$$

With the new solution encoding constraints (5)-(9) should be substituted with (13)-(15). This new solution encoding is used in both methods presented in the following sections. For given σ , and δ_{ij} distribution of the load can be obtained from formulation (1)-(4),(13)-(15), (10)-(12) simplified to a linear program.

3.2 Branch and Bound Algorithm

A branch and bound algorithm (BB) is a standard technique applied in solving combinatorial optimization problems. In BB algorithm a branching rule divides the set of possible solutions until distinguishing unique solutions. The pruning (or bounding) rule eliminates sets of solutions which are certainly not better than some already known solution, or are infeasible.

In our problem one has to determine a sequence of communications, and chunk overlapping. Communication sequences were built by appending a new processor to some already constructed leading sequence. For example, sequence $\sigma = (P_a, \dots, P_z)$ represents all the solutions beginning with communication sequence σ . This set is partitioned by appending a communication to any processor from set $\{P_1, \dots, P_m\}$. Thus, the set of solutions represented by σ is branched into subsets of solutions beginning with sequences: $(P_a, \dots, P_z, P_1), \dots, (P_a, \dots, P_z, P_m)$. For each communication sequence chunk overlapping on the used processors must be decided. All overlaps possible in the new encoding were enumerated in the following way. For processor P_i overlap is a vector $(\delta_{i1}, \dots, \delta_{in_i})$. A sequence $(\delta_{i1}, \dots, \delta_{ij})$ encoding the overlap for the first j chunks received by P_i , was branched into overlap encoding strings $(\delta_{i1}, \dots, \delta_{ij}, 0), \dots, (\delta_{i1}, \dots, \delta_{ij}, n_i - \max\{j, z_{ij}\})$. Thus, the BB algorithm uses a double branching scheme: for communication sequences, and for the chunk overlaps.

Enumeration of possible solutions was pruned by two methods. For a given sequence σ a lower bound $LB(\sigma)$ on the schedule length was calculated as follows. The startup times in σ were summed up: $\tau_1 = \sum_{i=1}^n S_{\sigma(i)}$. The maximum load V' that could be processed during the communication startup times is $V_1 = \sum_{i \in \sigma} (\tau_1 - \sum_{j=1}^{g(i)} S_{\sigma(j)}) / A_i$, where $g(i)$ is the index of the first communication to processor P_i in σ . The load must be sent from the originator in time at least $\tau_2 = V \min_{i=1}^m \{C_i\}$. In parallel with this communication, at most $V_2 = \sum_{i=1}^m \frac{\tau_2}{A_i}$ units of load could be processed. If $V_3 = V - V_1 - V_2 > 0$, then this remaining load V_3 will be processed in time at least $\tau_3 = \frac{V_3}{\sum_{i=1}^m \frac{1}{A_i}}$. The lower bound is equal to $LB(\sigma) = \tau_1 + \tau_2 + \max\{0, \tau_3\}$. Let T be the length of some already known solution. If $T \leq LB(\sigma)$ then successors of σ were discarded. Another mechanism used in sequence elimination was based on the maximum

memory $MEM(\sigma) = \sum_{i=1}^n B_{\sigma(i)}$ which could possibly become available in σ . If $MEM(\sigma) < V$ then, it means that memory available for holding the load is insufficient, and communication sequence is too short and must be expanded. In such a case the enumeration of the various overlap values was not attempted for the given σ . Observe that there are $O(m^n)$ communication sequences of length n for m processors, and for each processor the number of possible ways of overlapping the communication chunks is also exponential in n_i . Thus, due to the high computational complexity an upper bound n_{MAX} on length n of generated sequences was also imposed. Note that this was done to make the BB algorithm more usable, and it was not needed to properly define the algorithm. Consequently, due to the use of n_{MAX} not in all cases was BB able to deliver an optimum, or even a feasible solution.

3.3 Genetic Algorithm

The genetic algorithm (GA) is also one of the standard techniques used in solving combinatorial optimization problems. GA is a randomized algorithm using a set of operators transforming a population of solutions in the direction of improving quality. GA is defined by the way of encoding the solution, the set of genetic operators, stopping criteria, and several implementation-dependent tunable parameters.

In our implementation of GA solutions are encoded as pairs of strings. The first string is a sequence of processor indices encoding communication sequence σ . The second string O is encoding overlap of chunks. More precisely, $O(i)$ is the value of $\delta_{\sigma(i)j}$, where j is the number of load chunks sent to processor $P_{\sigma(i)}$ up to the i -th chunk sent off the originator. The lengths of σ, O are equal, and can be adjusted by GA to construct the best solution. Contents of strings σ, O is sufficient to formulate a linear program calculating α_i s as defined in Section 3.1. Fitness of the solution (called a chromosome) is measured as the value of schedule length T_{max} also obtained from the linear program for the given σ, O .

The genetic operators of GA applied here are selection, crossover, and mutation. The selection of the chromosomes for the new population is done by a combination of elitist and roulette wheel method and is strongly connected with the crossover operation. Each chromosome is selected with probability $\frac{1}{T_{max}^j} / \sum_{j=1}^G \frac{1}{T_{max}^j}$, where T_{max}^j denotes the schedule length for chromosome j , and G is the size of the population. The total number of selected parents is Gp_C , where p_C is a tunable algorithm parameter called crossover probability. In the crossover operation the selected parents are randomly paired and combined. For example, let $[(\sigma_1(1), \dots, \sigma_1(n')), (O_1(1), \dots, O_1(n'))]$ and $[(\sigma_2(1), \dots, \sigma_2(n'')), (O_2(1), \dots, O_2(n''))]$ be two parent solutions, with communication sequence lengths n', n'' , respectively. Let $k, l \leq m$ be two randomly chosen crossover points. The two offspring solutions are encoded in strings $[(\sigma_1(1), \dots, \sigma_1(k), \sigma_2(l), \dots, \sigma_2(n'')), (O_1(1), \dots, O_1(k), O_2(l), \dots, O_2(n''))]$, and $[(\sigma_2(1), \dots, \sigma_2(l), \sigma_1(k), \dots, \sigma_1(n')), (O_2(1), \dots, O_2(l), O_1(k), \dots, O_1(n'))]$. Note that because of choosing two crossover points l, k the offspring string lengths may be different than in their parents. The rest of the new population is se-

lected by elitist method so that the best $(1 - p_C)G$ chromosomes are always preserved. The elitist component in the selection was necessary because very often the difference in solution fitness is small, and the best solutions may be lost in the randomized selection based on the schedule length only.

Mutation changes $E(t)p_M$ random genes (i.e. pairs $(\sigma(i), O(i))$) in the population to different values. Here $E(t) = \sum_{j=1}^G n^j(t)$ is the total number of genes in the population in generation t , $n^j(t)$ is the length of chromosome j in iteration t , and p_M is a tunable algorithm parameter called mutation probability.

The algorithm stops after a fixed number of iterations it_1 . There is also a limit it_2 on number of iterations without an improvement in the quality of the best solution found so far. If iteration limit it_2 is reached before it_1 , then the population is replaced with randomly generated chromosomes and the search is started from scratch (the best solution found so far is recorded).

GA is a randomized algorithm which parameters must be tuned. We applied the following procedure. A set of 200 random instances with $m = 3, \dots, 5$, $V = 20$, B_i uniformly distributed in $[0, 10]$, A_i, C_i, S_i uniformly distributed in $[0, 1]$, were generated and solved to the optimum by BB. The average relative distance of the schedule length T_{max} from the optimum length was the measure of the tuning quality. The tunable parameters were selected one by one. The process of selecting the tunable parameters is illustrated in Fig.3. Intuitively, a big population size G should allow for finding good solutions in small number of iterations. On the other hand, maintaining big populations is computationally expensive. The populations size $G = 20$ was selected as a compromise between the speed of convergence to the near-optimum solutions, and the computational complexity (cf. Fig.3a). To select the crossover probability mutation operator was switched off. Crossover probability $p_C = 0.8$ was selected (Fig.3b). It turned out that majority of the population (80%) are offspring. Hence, it can be concluded that crossover is an effective optimization operator. After fixing G and p_C , mutation probability $p_M = 0.1$ was chosen (Fig.3c). In Fig.3d quality of tuning for various combinations of maximum number of iterations, and iterations without quality improvements are shown. Note that improving the average solution quality by 0.4% results in nearly 6-fold increase of the execution time. Hence, $it_1 = 100$, $it_2 = 10$ were selected as a compromise between quality and complexity.

3.4 BB vs GA comparison

Before proceeding to the further study of the features of the optimum and near-optimum solutions let us discuss advantages, and limitations of both solution algorithms.

BB guarantees optimum solutions, however, at considerable computational cost. In Fig.4 we compare average execution time of BB and GA. In the case of BB execution time is shown as function of n_{MAX} (Fig.4a). We use n_{MAX} because it turned out that the size of the search tree in BB is determined mainly by the limit n_{MAX} . The minimum communication sequence length

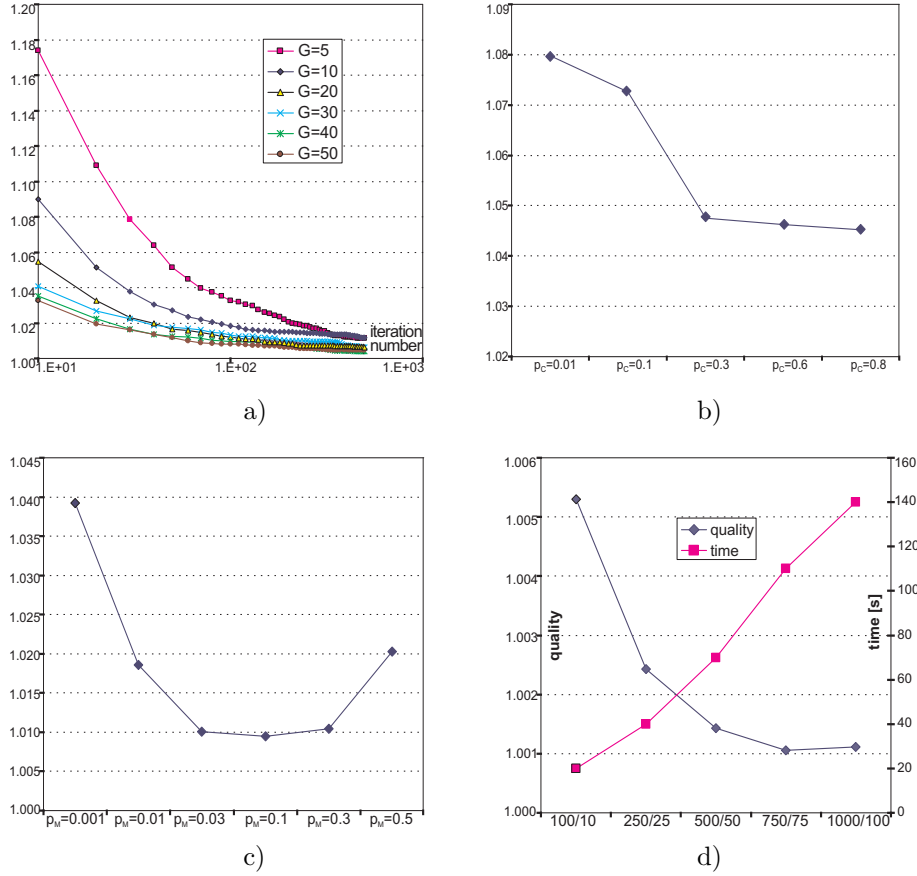


Figure 3: GA tuning. a) Solution quality vs. population size G , b) solution quality at 100th iteration vs. p_C , c) solution quality at 100th iteration vs. p_M , d) solution quality and execution time for various iteration limits it_1/it_2 .

$n_{MIN} = \frac{V}{\max_i \{B_i\}}$, is hardly ever the length of the best sequence, or the depth of the BB search tree. To be certain that the best communication sequence obtained in BB is indeed optimum it must have length at most $n_{MAX} - 1$. Instances satisfying this condition are easier to solve than the instances which force BB to search a tree as deep as n_{MAX} , and presenting the execution times in the function of the guaranteed optimum communication sequence length would not represent real execution time of BB. As it can be seen even average execution time of BB for $n_{MAX} = 7, m = 8$ is of order of one day. Hence, BB is not acceptable as a tool for studying features of great numbers of even moderate size instances. For the GA, execution time is shown vs. the length of the best obtained communication sequence. In Fig.4b execution time vs. the number of processors m is shown.

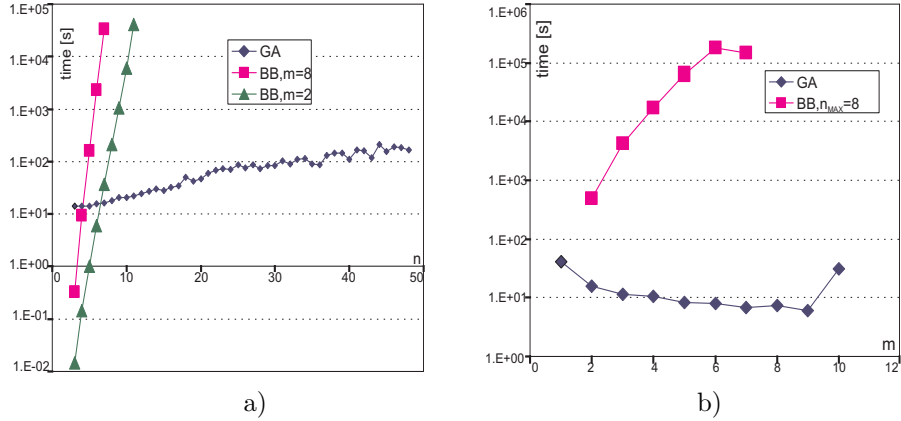


Figure 4: GA and BB execution times. a) vs. sequence length, b) vs. processor number m .

From the tuning process described in the previous section we conclude that GA is capable of delivering high quality solutions on average. As it can be seen in Fig.4 the running time of GA is much shorter than for BB. A disadvantage of GA as an analysis tool is that it is a randomized algorithm. In the limit of infinite iteration number, all feasible solutions are reachable in a process of random transformations of the solutions. The randomness in GA may have two disadvantages. Solutions which are not optimum may be too easy to find by GA, which may give wrong indications on the nature of the solved problem. On the other hand, solutions which have complex structure may be too improbable to be built in finite iteration number. For example, the communication sequence may include some processor which is not present in the optimum solution because the probability of selecting any processor in the sequence is high. Conversely, it is very unlikely that GA builds a long repetitive pattern of communications in the communication sequence because the probability of generating it decreases exponentially with the sequence length (more precisely it is $\frac{1}{m^n}$). Moreover, for the same instance GA may return different solutions in consecutive runs. For example, for a set of 45 random instances each solved 20 times, the quotient $\frac{T_{max}}{\overline{T_{max}}}$ where $\overline{T_{max}}$ is an average schedule length in all runs for a single instance, had the coefficient of variation was 6%, and average (over all quotients $\frac{T_{max}}{\overline{T_{max}}}$) was 0.9997.

We finish this section with a conclusion that BB is nearly unusable even on very moderate size instances. GA has much shorter execution time, and in the range in which it could be compared against BB, the quality of the GA solutions is very good. Though GA has its limitations it is the only tool at hand capable of solving bulk numbers of instances in reasonable time. Therefore, we will use GA as a replacement of BB in the analysis of the scheduling problem features.

4 Analysis of the Problem

In this section we present results of the investigation on the characteristic of the near-optimum solutions of our scheduling problem. We mainly concentrate on the features in the combinatorial part of the solutions: the communication sequence σ , and the vector of overlaps O . The studied features include:

- the need and the extent of the overlap,
- the length of communication sequence,
- number of used processors,
- the set of used processors,
- the chunk sizes,
- parameters of instances which make them easy, or hard, to solve.

We will draw conclusions both analytically and on the basis of experimental results. Both GA and BB in the simulations use `lp_solve` LP package [10]. Over 30000 instances were solved in the experiments on clusters of 15-75 PCs with Linux. Unless stated otherwise, the test data were generated in the following way. In the experiments involving analysis of the influence of system parameters A, B, C, S on solution characteristic, instance parameters A, B, C, S , were generated from $U(0, 1]$, i.e. uniform distribution with in range $(0, 1]$. The number of processors was generated from $U[1, 10]$, and all experiments were repeated for $V \in \{2, 5, 10, 20, 50\}$. In the experiments concerning certain parameter (say A) the parameter was fixed to a given value on all processors (e.g. $\forall i, A_i = 0.01$), and the remaining parameters were generated as described above. For each combination of V , and certain value of the parameter (e.g. $A_i = 0.01$) 1000 instances were generated. In the following sections we analyze the above features of the solutions.

4.1 Depth of Overlap

By the depth of overlap we mean the number of the load chunks which interfere with each other. The depth of overlap is expressed by the values of δ_{ij} which can be converted to the more convenient values of the span of overlap $z_{ij} - j + 1$ for each chunk j on processor P_i , where z_{ij} (see Section 3.1) is the number of the last chunk overlapping with chunk j on P_i . The existence of the overlap means that load must accumulate on the processors. It is of practical importance to verify if the accumulation of the load is actually necessary, and to what degree.

4.1.1 Single processor considerations

Here we analyze the case of one computing processor ($m = 1$). Despite simplicity, this problem is not trivial because to construct a schedule one has to decide on the overlap of the load chunks, and their sizes.

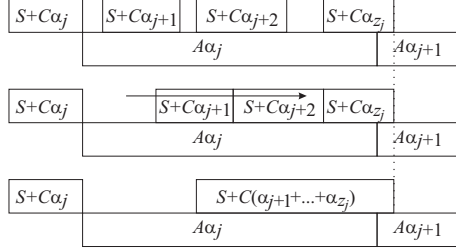


Figure 5: Merging overlapping chunks in Observation 2.

Let us start with several simple observations which can help in reasoning about our problem. For simplicity presentation we drop the subscripts related to processor indices. We denote processor parameters as A, B, C, S , and the overlap of chunk j by δ_j . It can be shown [12] that if there are no memory limitations, then there are no idle times in the computations nor in the communications (also for $m > 1$). It is not the case in our problem.

Observation 1 *There are optimum sequences with idle times in the computation and in the communication.*

Proof. Suppose $A = 1, B = \frac{V}{2}, C = 0, S = M$, where $M \gg V$ is a big constant. The minimum number of communications is $n_{MIN} = \frac{V}{B} = 2$, for which schedule length is $T_{max} = 2M + V$. There is an idle interval of length $\frac{V}{2}$ in the communications, and two idle intervals of length M in the computations. Idle times in computation cannot be closed because any load which feasibly fits in memory size $B = \frac{V}{2}$ is processed in shorter time than the startup time $S = M \gg V$. Suppose that we want to close the idle interval in the communications. This requires that two consecutive load chunks fit in memory of size $B = \frac{V}{2}$, and at least three chunks must be sent from the originator. Then schedule length is at least $3M$. Since M can be arbitrarily big in relation to V , the difference between the length of a schedule with idle times and the length of the schedule without idle times in communication can be arbitrarily big in absolute terms. \square .

We will be saying that solutions for which chunks overlapping by not more than 1, i.e. $\forall j, \delta_j \leq 1$, have overlap at most 1. If $\forall j, \delta_j = 1$, then we will say that a solution has overlap 1.

Observation 2 *On one processor there is no need for more than overlap 1.*

Proof. We will prove that there are optimum solutions for which $\delta_j \leq 1$. Suppose there is an optimum solution with overlap greater than one and j is the first chunk for which $\exists j, \delta_j > 1$. Chunk $z_j > j + 1$ is the last chunk overlapping with chunk j . Consider a constructive transformation (see Fig.5). First communications of the chunks $j + 1, \dots, z_j - 1$ are shifted right (later) until they are sent without idle time just before chunk z_j . Second, include the load

of chunks $j+1, \dots, z_j-1$ in the message of z_j , and drop chunks $j+1, \dots, z_j-1$ altogether. This operation is feasible because the optimum solution must satisfy condition (15) demanding that all chunks j, \dots, z_j fit in memory. Thus the load fits in the memory no matter whether it is sent in chunks $j+1, \dots, z_j-1$, or just in chunk z_j . This operation does not change schedule length because chunks $j+1, \dots, z_j-1$ overlap with j , and hence, are not being processed when chunk j is computed. The overlap δ_j has been reduced to at most 1. This operation may be repeated iteratively to all chunks to remove any overlap greater than 1. \square .

Observation 3 *The schedule for sequence with the n_{MIN} can be at most twice as long as the optimum schedule.*

Proof. Schedule length T_{max} for a sequence with smallest $n_{MIN} = \lceil \frac{V}{B} \rceil$ is not greater than $n_{MIN}S + CV + AV$. On the other hand, for the optimum solution, $T_{max}^* \geq n_{MIN}S + CV$ and $T_{max}^* \geq AV$.

1. If $AV \leq n_{MIN}S + CV$ then $\frac{T_{max}}{T_{max}^*} \leq \frac{n_{MIN}S + CV + AV}{n_{MIN}S + CV} \leq 2$.
2. If $AV \geq n_{MIN}S + CV$ then $\frac{T_{max}}{T_{max}^*} \leq \frac{n_{MIN}S + CV + AV}{AV} \leq 2$. \square .

Observation 4 *In solutions with overlap 1 chunk sizes may be coupled.*

Below we explain what we mean by *coupling* of chunk sizes. A schedule with overlap 1 has $\forall j, \delta_j = 1$, and the chunks overlap with their direct predecessor and direct successor (if any). If chunk 1 has size α_1 , then by (15) chunk 2 has size at most $\alpha_2 \leq B - \alpha_1$, chunk 3 has size at most $\alpha_3 \leq B - \alpha_2 = \alpha_1$, etc. Thus, if chunks have all their maximum sizes, then the size of all chunks is in fact determined by a single variable α_1 . The size of processed load is $\frac{n}{2}B$ if communication sequence has even number n of messages, or it is $\frac{n}{2}B + \alpha_1$ if n is odd. Chunk sizes are coupled if $n = \lceil \frac{2V}{B} \rceil$, and the overlap is 1 (thus number of chunks is chosen tightly for the given overlap). We will say that solutions with $n = \lceil \frac{2V}{B} \rceil$, overlap 1, for $m = 1$ are coupled.

Observation 5 *The coupled solutions are not arbitrarily bad.*

Proof. For the optimum solution $T_{max}^* \geq \lceil \frac{V}{B} \rceil S + CV = n_{MIN}S + CV$. For a coupled solution we have $T_{max} \leq \lceil \frac{2V}{B} \rceil S + CV + AV = n_2S + CV + AV$ and $n_2 \leq 2n_{MIN}$.

1. If $AV \leq CV$ then

$$\frac{T_{max}}{T_{max}^*} \leq \frac{n_2S + CV + AV}{n_{MIN}S + CV} \leq 1 + \frac{n_{MIN}S + AV}{n_{MIN}S + CV} \leq 2$$

2. If $n_2S + CV \leq AV$ then

$$\frac{T_{max}}{T_{max}^*} \leq \frac{n_2S + CV + AV}{AV} \leq 1 + \frac{n_2S + CV}{AV} \leq 2$$

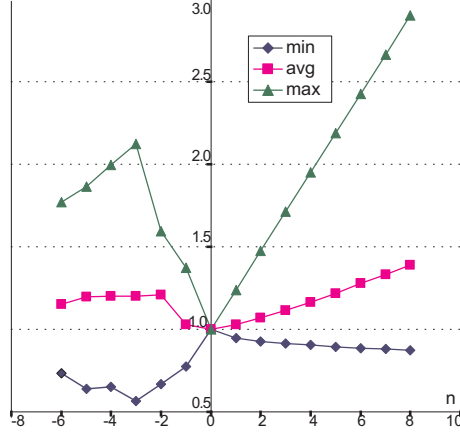


Figure 6: $m = 1$, quality of the solutions with various communication sequence lengths, the best overlap, relative to coupled solutions.

3. If $CV \leq AV \leq n_2S + CV$ then

$$\begin{aligned} \frac{T_{max}}{T_{max}^*} &\leq \frac{n_2S + CV + AV}{n_{MIN}S + CV} \leq \frac{2(n_2S + CV)}{n_{MIN}S + CV} \leq \frac{2(2n_{MIN}S + CV)}{n_{MIN}S + CV} \leq \\ &\leq \frac{4(n_{MIN}S + CV)}{n_{MIN}S + CV} = 4 \end{aligned}$$

Let us note that if the number of communications is not greater than $k \times n_{MIN}$, then $\frac{T_{max}}{T_{max}^*} \leq 2k$. \square .

The above observation gives an indication on the quality of schedules with $n = \lceil \frac{2V}{B} \rceil$ and overlap 1. In Fig.6 quality of schedules for $m = 1$, various sequence lengths, and the best overlap chosen by BB algorithm is shown. The coupled solution quality is used as a reference, and are represented by the points at the coordinates (0, 1). For example, shorter sequences are shown on the negative part of horizontal axis. The best, the worst, and an average distance from the coupled solution is shown. The results in Fig.6 represent 888 randomly generated instances with $A, C, S \sim U[0, 1]$, $B \sim U(0, 10)$, $V = 10$. As it can be seen, typically the best solutions are not very much better than the coupled ones. Increasing n beyond $\lceil \frac{2V}{B} \rceil$ is not reducing schedule length more than by approx. 13%. Thus, on average coupled solutions provide a simple and efficient method of solving the combinatorial part of our problem on $m = 1$ processor. Let us observe that optimum communication sequence length n can be smaller or greater than $\lceil \frac{2V}{B} \rceil$ depending on the instance.

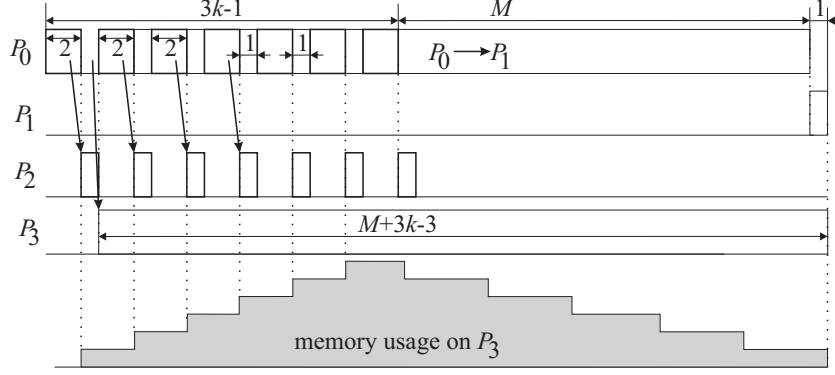


Figure 7: An instance with arbitrary overlap in Observation 6.

4.1.2 Overlap on Multiprocessors

In the preceding section it has been observed that no overlap greater than 1 is necessary if $m = 1$. However, for $m > 1$ arbitrarily big overlap may be necessary in optimum solutions.

Observation 6 *There are optimum solutions with arbitrarily big overlap.*

Proof. Let k, M be two integers, where $k > 5$ is even, and $2^k - 3k - 1 > M > 3k + 1$. Consider the following example: $m = 3, V = 2^{2k} + k2^k + 1, A_1 = \frac{1}{2^{2k}}, B_1 = 2^{2k}, C_1 = 0, S_1 = M, A_2 = \frac{1}{2^k}, B_2 = 2^k, C_2 = 0, S_2 = 2, A_3 = M + 3k - 3, B_3 \geq V, C_3 = k - 1, S_3 = 0$. We want to build a schedule of length $T \leq M + 3k$. We will show that also no shorter schedule may exist.

To process V one activation of P_1 is necessary. P_1 cannot be substituted by any other processor because only P_1 has sufficient memory buffer, and speed to process load 2^{2k} in time shorter than T . Indeed, load 2^{2k} is computed on P_1 in time 1. P_1 cannot be activated more than once because $2S_1 = 2M > T$. If P_2 were to substitute P_1 then it would require at least $\frac{B_1}{B_2} = 2^k$ communications of length at least $\frac{B_1 S_2}{B_2} = 2^k > T$.

We will prove now that processor P_2 must receive k messages processing of which is not overlapping.

Consider the load $V - B_1 = k2^k + 1$ remaining to be processed by P_2, P_3 . In time T processor P_3 is capable of processing at most $\frac{T}{A_3} = \frac{M+3k}{M+3k-3} = 1 + \frac{3}{M+3k-3} < 2 < 2^k$ units of load. Note that $\frac{3}{M+3k-3}$ can be arbitrarily small by the selection of M, k . Thus, to process the remaining load P_2 must receive at least k messages. If in the k messages each one carries load B_2 than the whole communication and computation can be feasibly performed in time $3k$ as shown in Fig.7.

On the other hand assume that P_2 receives more than k messages. Suppose that processing of none of the load chunks overlaps. Then we have communication time at least $(k+1)S_2 = 2k+2$. At least $k2^k - \frac{3}{M+3k-3}$ units of load must

be processed by P_2 . Computation of at most B_2 load units can be performed in parallel with startup S_1 . Hence, computation and communications of P_2 together with the startup time S_1 take at least time $S_1 + (k+1)S_2 + A_2(k2^k - \frac{3}{M+3k-3} - B_2) = M + 2k + 2 + k - 1 - \frac{3}{2^k(M+3k-3)} > M + 3k$. Thus, if the number of messages received by P_2 is greater than k , then processing of load chunks on P_2 must overlap.

We will calculate now the length of the communication time when the chunks sent to P_2 overlap. Consider some sequence of load chunks which processing is overlapping. Let the first load chunk i overlap with the next δ_i chunks. To obey memory limits the group of δ_i consecutive chunks may contain load at most B_2 . The next group of chunks starting with chunk $\delta_i + 1$ is independent in this sense that they may carry another volume of at most B_2 . Let l_a be the number of groups of chunks overlapping by a . Let Δ be the greatest overlap. A group of a overlapping chunks requires communication time at least $S_2 a$. If some chunk i is not overlapping (denoted overlap $\delta_i = 0$) with any other chunk then its communication and computations (not overlapped by other communication to P_2) last $S_2 + A_2 \alpha_i = S_2 + \frac{\alpha_i}{2^k}$. In the following discussion it will be convenient to assume that each chunk with overlap 0 transfers load of size B_2 . We explain that this transformation does not increase schedule length. Let V' be the amount of load that should be added to the chunks with overlap 0 to fill the buffer completely, i.e. $V' = \sum_{\{i:\delta_i=0\}} (B_2 - \alpha_i)$. This missing amount of load must be processed in chunks with the overlap greater than 0. If we shift load V' from chunks with overlap 0 to the chunks with overlap greater than 0, then we shorten the idle time in communications to P_2 by $V' A_2 = \frac{V'}{2^k}$. Simultaneously, we increase the number of communications by $\lceil \frac{V'}{B_2} \rceil$, which take time $S_2 \lceil \frac{V'}{B_2} \rceil = 2 \lceil \frac{V'}{2^k} \rceil$. Thus shifting the load from chunks with overlap 0 to other chunks is not reducing the duration of communications. Hence, to calculate a lower bound on the schedule length we may assume that chunks with overlap 0 contain load B_2 . To process on P_2 the required amount of load the number of independent groups of overlapping chunks must be at least k as shown above, i.e. $\sum_{\delta=0}^{\Delta} l_{\delta} = k$. The total communication time is at least

$$\begin{aligned} & (S_2 + 1)l_0 + S_2 \sum_{\delta=1}^{\Delta} (\delta + 1)l_{\delta} = \\ & (S_2 + 1)l_0 + \sum_{\delta=1}^{\Delta} ((S_2 \delta - 1)l_{\delta}) + \sum_{\delta=1}^{\Delta} ((S_2 + 1)l_{\delta}) = \\ & (S_2 + 1) \sum_{\delta=0}^{\Delta} l_{\delta} + \sum_{\delta=1}^{\Delta} ((S_2 \delta - 1)l_{\delta}) = \\ & (S_2 + 1)k + \sum_{\delta=1}^{\Delta} ((S_2 \delta - 1)l_{\delta}) = \\ & 3k + \sum_{\delta=1}^{\Delta} ((2\delta - 1)l_{\delta}). \end{aligned}$$

Note that $\sum_{\delta=1}^{\Delta} ((2\delta - 1)l_{\delta}) > 0$. If $\exists \delta > 1, l_{\delta} > 0$ then the communications to P_2 are longer than $3k$, and hence it is infeasible to send enough load to P_1 , and P_2 . Therefore, P_2 can receive at most k messages processing of which must not overlap.

Since P_2 receives k messages at most $k2^k$ units of load are processed on P_2 . Processor P_3 must compute the remaining amount of $V - B_1 - k2^k = 1$ units of load. Transferring the k load chunks to P_2 takes communication time $S_2 k = 2k$. Since the chunks on P_2 are not overlapping, at $k - 1$ intervals of length $B_2 A_2 = 1$

Table 2: Relative frequency of the overlaps in all the chunks

overlap	0	1	2	> 2
frequency	0.835	0.154	0.010	0.001

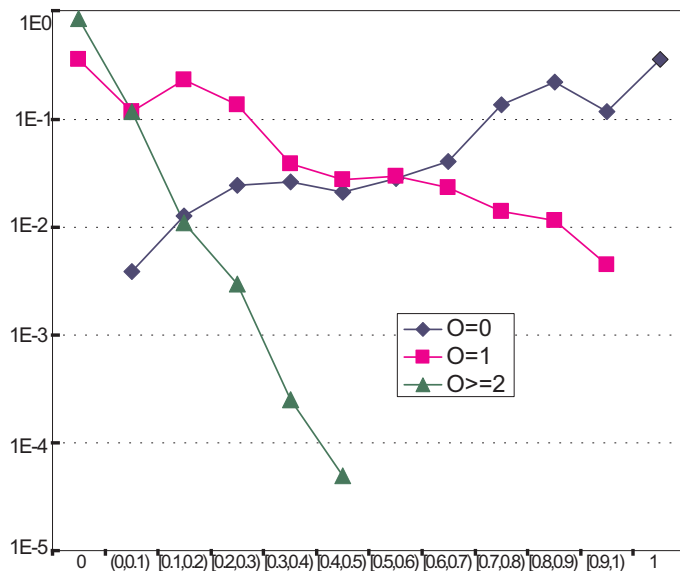


Figure 8: Histogram of overlap frequencies in the instances.

are free for communications to P_3 . The communications to P_1 and P_2 together with the necessary idle times in communication to P_2 take time $M + 3k - 1$. Considering the fact that the load on P_1 or on P_2 must be computed in at least 1 unit of time, P_3 cannot receive its load before P_1 , or P_2 . Thus, load for P_3 will be delivered in at most k chunks. Moreover the schedule cannot be shorter than $M + 3k + 1$. Communication time of each load chunk for P_3 is not longer than 1, therefore load of size at most $\frac{1}{k-1}$ can be transferred to P_3 in a single message. Computing 1 unit of load on P_3 lasts $M + 3k - 3$. Thus there can be no more than three units of idle time on P_3 . In parallel with startup time S_1 processor P_3 must compute at least load $\frac{M}{M+3k-3}$. Since messages convey at most load $\frac{1}{k-1}$, the number of messages waiting to be processed on P_3 when communication to P_1 starts must be at least $\frac{(k-1)M}{M+3k-3}$ which tends to $k - 1$ as M tends to infinity. We conclude that it is possible to construct an optimum schedule which requires arbitrarily deep overlap. \square

From the above observation we conclude that the overlap can be arbitrarily deep. Results collected from solutions by GA of 19953 randomly generated instances with $A, B, C, S \sim U[0, 1]$, $m \sim U[1, 10]$ and $V \in \{2, 5, 10, 20, 50\}$

indicate, however, that the depth of the overlap is not very big. Table 2 lists the relative depth of the overlap of all chunks in all sequences of the solutions generated by GA for the above instances. A more detailed view of the chunk overlaps is presented in Fig.8. The vertical axis is the relative frequency of instances with certain fraction of chunks of with certain overlap. For example (see the rightmost box "1" for overlap $O = 0$), for chunks with overlap 0, approximately 36% of all instances have only chunks with overlap 1. The absence of a point in box "0" for overlap 0 means that there were no instances without a chunk with overlap 0. The number of solutions for which the chunks with overlap 1 are 90% to 99.99% of all the chunks in the communication sequence is approx. 0.4% of all instances (box "[0.9,1]" for overlap $O = 1$). On the other end, approx. 36% of all instances have no chunk with overlap 1. Finally, overlaps 2 and bigger are very rare: approx. 87% solutions have no chunk with overlap 2 or greater, and only 0.005% instances have solutions with the number of chunks with overlap at least 2, in more than 40% of all the chunks in the solution. For the above experimental results we can conclude that overlap deeper than 1 is rare, because it constitutes approx. 1% of all chunks in all solutions.

The analysis of the depth of the overlap leads to the following conclusions: On a single processor the overlap is in $\{0, 1\}$, the solutions with overlap 0, or 1 cannot be arbitrarily bad, and a solution with $n = \lceil \frac{V}{B} \rceil$, and $\forall j, \delta_j = 1$ is good on average. For multiple processors ($m > 1$), the the overlap may be arbitrarily deep, but overlaps greater than 1 are rare in practice. Moreover, results obtained for $m = 1$ cannot be transferred to the multiprocessor case.

4.2 Length of the communication sequence

One of the important characteristics of the solution is the number of communications n . The length of the communication sequence depends on V , and B_i s. Therefore, it seems reasonable to use some reference number of communications. Let us assume that the reference number of communication is n_{MIN} . We start with some observations.

Observation 7 *A communication sequence with minimum number of chunks $n_{MIN} = \frac{V}{\max_i \{B_i\}}$ can be arbitrarily bad for schedule length.*

Proof. Consider an example: $m = 2, S_1 = 1, C_1 = 0, A_1 = 1, B_1 = 1, S_2 = M, C_2 = 0, A_2 = 1, B_2 = V$, where M is a big constant. The minimum number of communications is $n_{MIN} = 1$, and it results in a schedule of length $M + 1$. On the other hand if P_2 is used only, then schedule length is $1 + V$, and the number of communications is $n = \lceil V \rceil$. The ratio of the two schedule lengths is $\frac{M+1}{V+1}$ which can be made arbitrarily big by selection of M , and V . \square .

Observation 8 *The length n of the optimum communication sequence can be arbitrarily big in relation to n_{MIN} .*

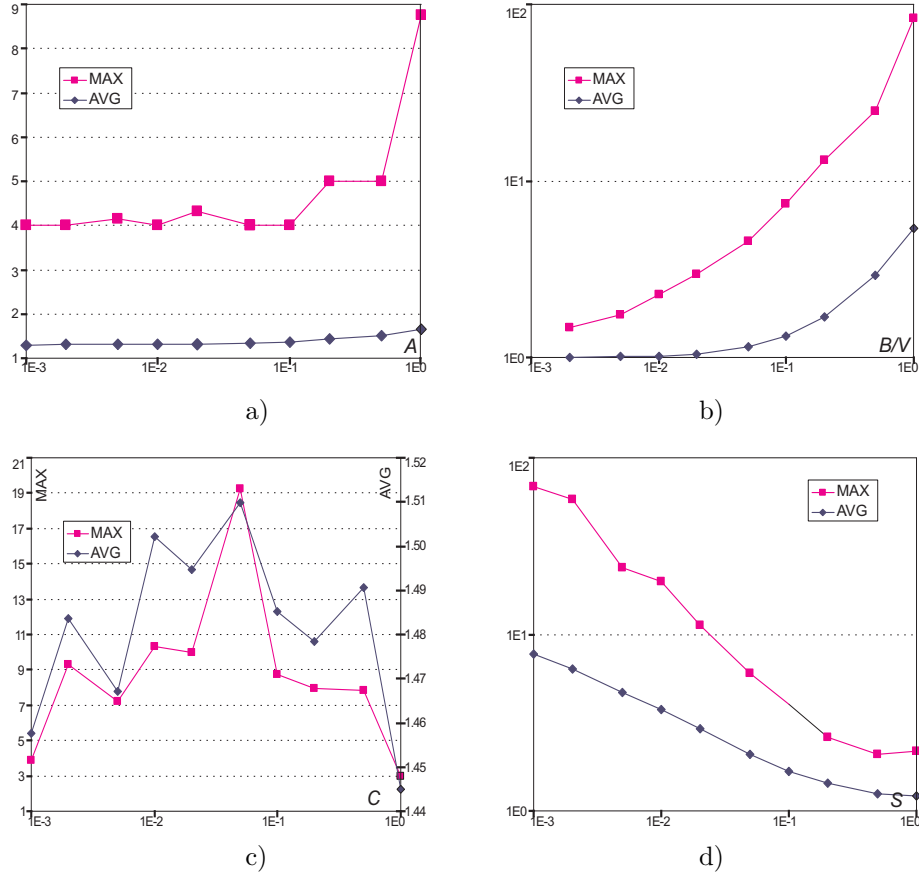


Figure 9: Relative sequence length $\frac{n}{n_{MIN}}$ in the solutions of GA a) vs. A , b) vs. $\frac{B}{V}$, c) vs. C , d) vs. S .

Proof. Note that the number of messages in the optimum communication sequence proposed in the previous proof can be arbitrarily big. \square .

Let us now analyze the length of communication sequences generated by GA. The values of relative communication sequence lengths $\frac{n}{n_{MIN}}$ are shown in Fig.9. In Fig.9a the average (AVG), and the longest (relative) communication lengths are shown for various A values. It can be seen that typically $\frac{n}{n_{MIN}}$ is not very big. On average $n \approx 1.39n_{MIN}$, which is calculated over all instances of changing A . The length of the sequence grows with A , which is especially evident for the biggest registered relative lengths. This phenomenon can be attributed to the way of calculating n_{MIN} . For example, for $V = 2$, and $B_i \in (0, 1]$ the expected n_{MIN} is 4, and in extreme cases it can be just $n_{MIN} = 2$. On the other hand, as processors get slower (A is increasing) it gets more and more profitable to use all m available processors. Thus, $\frac{n}{n_{MIN}}$ grows with A . This increase is

stronger for small V , and weaker for bigger V . In Fig.9b similar dependence is shown for changing $\frac{B}{V}$. The length of the communication sequence quickly increases with $\frac{B}{V}$. It is because, on one hand, for $\frac{B}{V}$ approaching 1, n_{MIN} is also approaching 1, but as in Observation 8 other parameters of the system make it profitable to build sequences with $n \gg 1$. On the other hand, as $\frac{B}{V}$ approaches 0, more and more short communications must be made to send the load off the originator. Each message carries cost of startup S . Therefore, communication costs and startup in particular, dominate in the schedule length. To minimize this cost it is advantageous to send as few messages as possible. Hence, n tends to n_{MIN} when $\frac{B}{V}$ is decreasing. Similar observations can be made for big values of S (cf. Fig.9d). For big S_i s it is profitable to send as few messages as possible. This, in turn, exposes the need for big communication buffers. The behavior of $\frac{n}{n_{MIN}}$ for small S must be contrasted with Fig.9a. When $S \approx \frac{1}{2}$ on average, as in Fig.9a, then $\frac{n}{n_{MIN}} \approx 1.39$. If $S = 0.001$, as in Fig.9d, then $\frac{n}{n_{MIN}} \approx 8$. This means that big startup time is a considerable disincentive to building long communication sequences. In Fig.9c dependence of $\frac{n}{n_{MIN}}$ on C is shown. Note that this figure has two vertical axes. The shapes of MAX , and AVG are similar, but for the average case the changes are in the range of approximately 5%. This should be surprising because multi-installment divisible load processing was introduced to reduce the time of initial waiting for load. Growing value of C should be an incentive to build shorter messages and longer communication sequences. This tendency can be seen only for small values of C . Yet, in our setting of the experiments expected value of the startup times is $\frac{1}{2}$ which is a disincentive to build long communication sequences as explained on the example of Fig.9a, and Fig.9d. Hence, the dependence of average $\frac{n}{n_{MIN}}$ on C is very weak. Moreover, with growing C the algorithm tends to compensate increasing communication costs by sending fewer messages. Thus, initial waiting for the load is meaningless compared to the whole communication cost.

From the above analysis of the communication sequence length we draw the following conclusions: Startup times S_i are important element of communication time, and they constitute main disincentive to build long communication sequences. For startup times of the same order as communication time per unit of load (C), or computation time per unit of load (A) communication sequences have lengths $\approx 1.4n_{MIN}$. For small S the sequences can be approximately 8-10 times longer on average than n_{MIN} . Moreover, S_i , and B_i are in a sense coupled in determining system performance: small B_i s expose costs of communications including startups, big S_i s expose the need for processors with big communication buffers.

4.3 Number of used processors

In this section we study the number m' of different processors used. This feature of a solution has a practical meaning. Considering big pools of processors available in contemporary grid and cluster systems it is of practical importance to know how many processors should be used, and how to adjust their num-

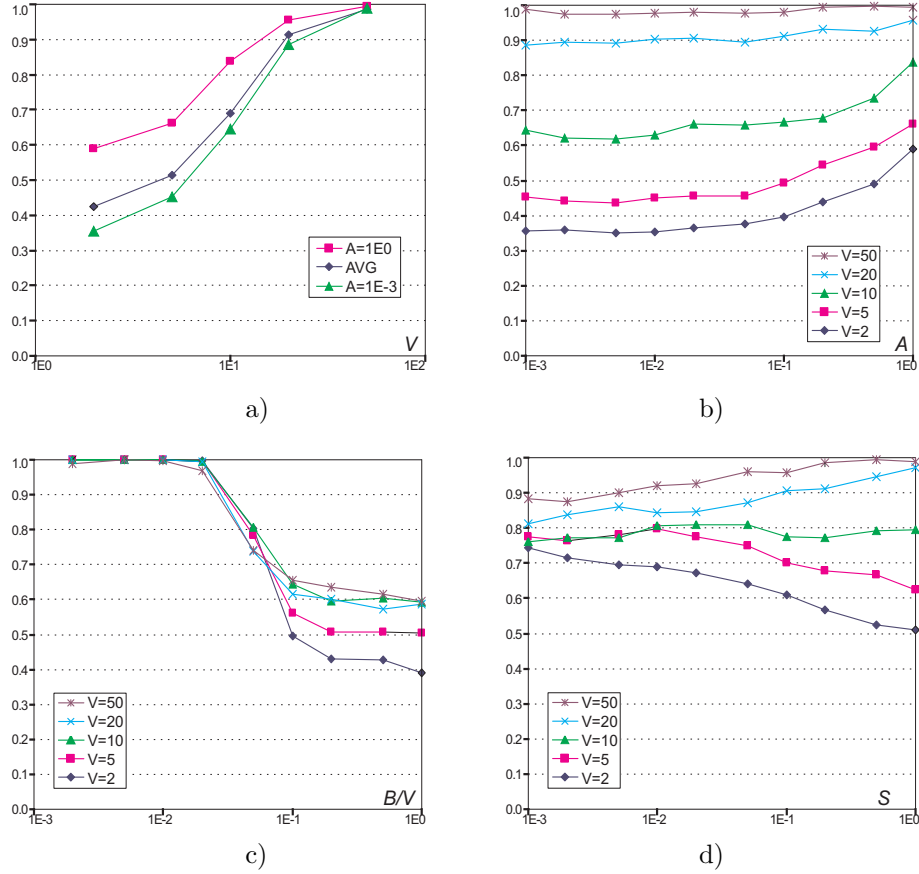


Figure 10: Relative number $\frac{m'}{m}$ of different used processor in the solutions of GA a) vs. V , b) vs. A c) vs. $\frac{B}{V}$, d) vs. S .

ber to changing characteristic of the system and application. Let us start with some observations. It is not difficult to coin instances where only one processor may be used (e.g. because all other processors have their startup times greater than schedule length) or all processors must be used (e.g. identical processors, no startup time, huge V). It is known fact from divisible load theory that if $\forall i, S_i = 0$, then computations can be started on arbitrary number of processors. On the other hand if $\exists i, S_i > 0$, then for single installment processing using all processors is a matter of sufficiently big volume of load V . Thus, it may be intuitively expected that the number of used different processors m' should grow with decreasing startups and increasing V .

Let us now analyze the features of GA solutions. Relations between the ratio $\frac{m'}{m}$ and selected parameters are shown in Fig.10. It can be seen in Fig.10a that with growing amount of load V the number of different used processors

is growing as could be intuitively expected. This result was confirmed in all experiments we performed. This has a practical consequence, that for bigger problems it is profitable to use more processors instead of relying on bigger number of load chunks n only.

The dependence of m' on A is show in Fig.10b. Only for small problem sizes (small V) does m' increase with A . For small V only a few chunks need to be sent. Therefore, for small A the algorithm minimizes schedule by selecting only a few processors with big memory buffers and fast communication link. For big A computing time dominates schedule length, and it is profitable to distribute and parallelize computations. Hence $\frac{m'}{m}$ is growing. On the other hand, for big V the number of chunks must be big anyway, computation time (mainly startup times S_i) is dominating over communication time, and A is less important in determining schedule length. Therefore, A is not influencing m' for big V .

In Fig.10c dependence of $\frac{m'}{m}$ on $\frac{B}{V}$ is shown. In our method of test instance generation average number of processors is ≈ 5 . Hence, for $\frac{B}{V} < \frac{1}{5}$ the memory space necessary to process load V is created by using many load chunks, and many processors working in parallel. On the other hand, when $\frac{B}{V} > \frac{1}{5}$ the size of memory is sufficient to process the whole load in just one installment. Therefore, good solutions tend to use only a few processors with fast communication and computation.

Fig.10d shows relation between S , V , and $\frac{m'}{m}$. With growing amount of load V the number of different used processors is increasing as in other experiments. For small V the number of different used processors decreases with S which is in accord with our earlier expectations. However, for big V the increasing S results in increased m' . This counterintuitive behavior partially can be explained by the way of generating test instances. Note that startup times of all processors are equal in the experiments depicted in Fig.10d. When V is big then the number of sent chunks must be also big. With growing S startup times dominate in the schedule length and other parameters, by which the processors differ, are becoming meaningless. Therefore, GA becomes myopic to the differences in processor parameters, and hence more processors are drawn to the solutions.

Dependence of $\frac{m'}{m}$ on C (not shown here) is very weak. This is a very surprising situation because in many DLT papers communication rate C was considered crucial for system performance. Only for small V and big C (close to 1) is the number of used processors slightly decreasing with growing C . This is a result of the startup time domination in communication time. Only for small V the number of messages is small and hence startups is small. Then, GA optimizes the schedule by using few efficient processors. This result does not eliminate C as an important performance determinant, as will be shown in the following study.

We finish this section with the following conclusions. The number of different used processors differs depending on the settings. In general it is increasing with V . In our experiment setting startup times dominated schedule length, especially when the number of chunks had to be big because V was big or B was small. When A is big and computation time is at least comparable with the

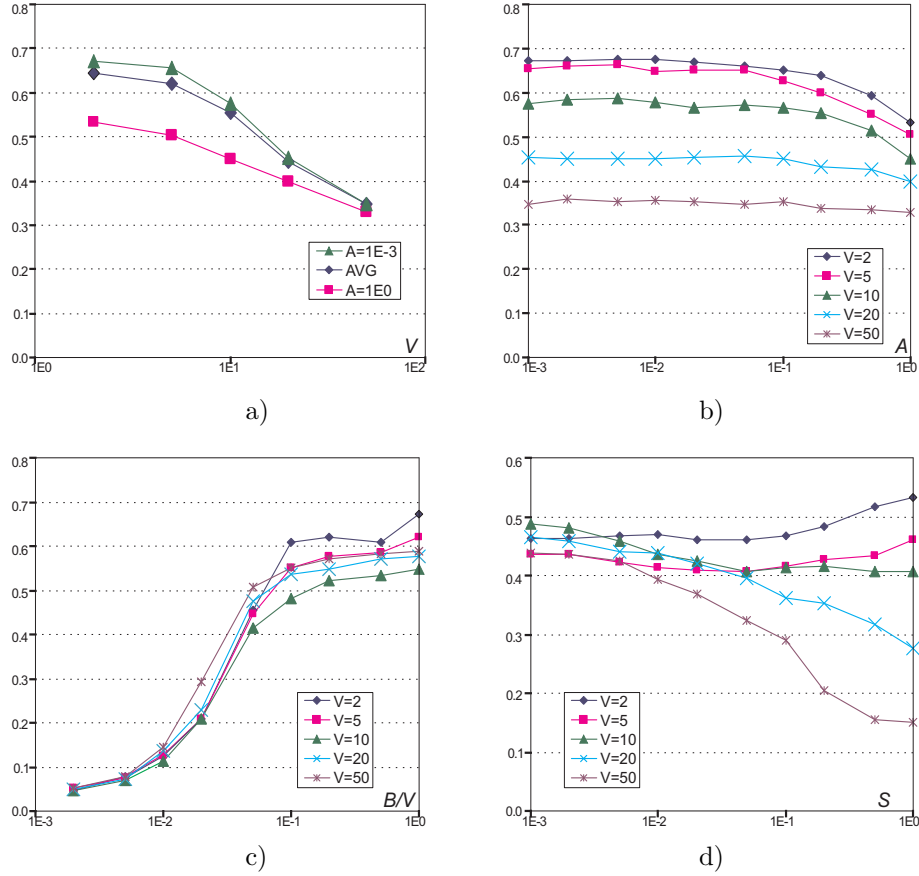


Figure 11: Gini index of the GA solutions a) GiL vs. V , b) GiL vs. A , c) GiL vs. $\frac{B}{V}$, d) $Gi\#$ vs. S .

communication time, then it is profitable to use many processors to parallelize computations. When C is big and its contribution to the communication time is comparable, or greater, than the contribution of the startup times, then it is profitable to choose few processors with small C .

4.4 Dominating Set of Processors

Observe that previous section considered only the number of different used processors, not the degree of participation in computations. Here we analyze distribution of the load between the processors. We want to determine if there is any inequality in the load distribution, and if it is the case, then what kind of processors dominate in the computations.

As the first tool in analyzing inequality in the load distribution we applied Gini index. Gini index is an indicator of some parameter deviation from the

uniform (flat) distribution. It is commonly used in economics to quantify inequality in wealth distribution. The closer Gini index is to 0 the more equal and uniform distribution of the load is. The closer Gini index is to 1, the more unequal distribution of the parameter is. We calculated Gini indices for the load received by the processors (which we will denote GiL), and for the number of communications (which we will refer to as $Gi\#$). For example, $GL = 1$ implies that whole load V is processed by a single processor. Selected results are presented in Fig.11. The general observation was that GiL , and $Gi\#$ have similar shape and demonstrate the same tendencies.

It can be seen in Fig.11a that GiL is decreasing with V which means that with growing size of the load the distribution of the load is becoming more and more equal. This situation has been observed in all experiments. Dependence of GiL on A is shown in Fig.11b. Only for small V does A influence load distribution. For small V the number of used processors is small and it is profitable to select the best of them, while for big V the number of load chunks must be big anyway which means that communication time is long and computation time (hence A) has little influence on schedule length. Consequently, for big V values of GiL do not depend on A . This situation is similar to Fig.10b depicting $\frac{m'}{m}$ vs A . A strong change of GiL with $\frac{B}{V}$ is observed in Fig.11c when $\frac{B}{V} \approx \frac{1}{5}$. For smaller values of $\frac{B}{V}$ the load distribution is more equal, for bigger $\frac{B}{V}$ the load distribution is more unequal. This means that for $\frac{B}{V} > \frac{1}{5}$ only one installment is sufficient to process the whole load. These results conform with the results depicted in Fig.10c. In Fig.11d $Gi\#$ is shown for changing S . Again, similarly to Fig.10d, with growing S the diversity of used processor sets depend on V . For small problem sizes it is profitable to use fewer processors, hence $Gi\#$ is big which signifies inequality. For big V the number of used processors is big, communication startup times dominate in the schedule length, and the algorithm becomes myopic to the differences in processor parameters, hence more of them get to the communication sequence, and the messages are distributed more equally.

Unfortunately, Gini index is hard to interpret. For example, it is hard to say if certain value of GiL , $Gi\#$ already represents inequality or not. Only general tendencies of changing inequality can be observed. Here, the tendencies of Gini index confirm the analysis of the number of used processors. Therefore, we applied one more index of load distribution inequality.

The second measure of processor domination in computations is based on the set of processors most frequently receiving the load or messages. Let p_{max} be the greatest number of load chunks sent to a single processor. Processors receiving at least $\frac{p_{max}}{2}$ messages will be called *communication frequent*. We will say that the set of all communication frequent processors is *communication frequent set*. Similarly, let V_{max} be the greatest total load received by any processor. We will say that some set of processors is *load frequent* if it includes all processors which receive at least $\frac{V_{max}}{2}$ units of load. The processors in load frequent set are called load frequent.

We want to examine how much load, and messages are sent to frequent

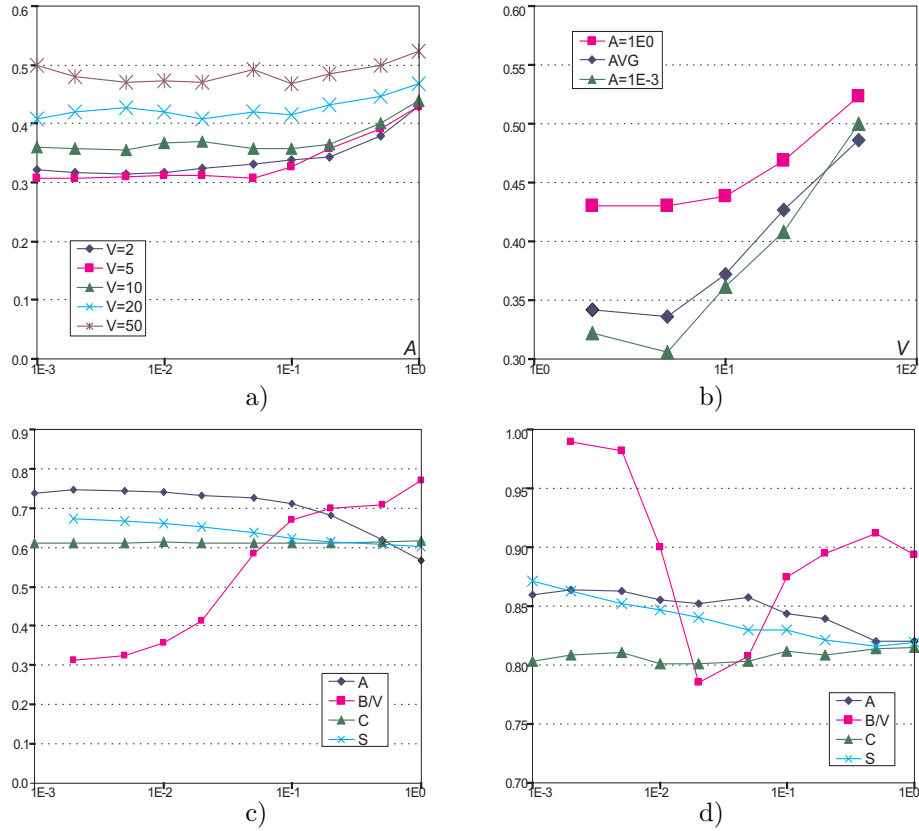


Figure 12: Load frequent processor sets in GA solutions. a) number of frequent processors vs. A , b) number of frequent processors vs. V , c) load of the most loaded processor and d) load of all the frequent processors vs. A , $\frac{B}{V}$, C , S .

processor sets. The results of this study are shown in Fig.12. All values shown in this figure are relative. Thus, processor numbers are shown with respect to m , and the loads are shown relative to V . In Fig.12c,d, the horizontal axes represent all parameters A , $\frac{B}{V}$, C , S , in range $[0, 1]$ for four different relations. A general observation is that the characteristics of *load frequent* sets are very similar for the *communication frequent* sets. Therefore, we will present the results for load frequent processors only and we will refer to them as to just frequent processors. Another general observation is that the functions of the number of load frequent processors in A (Fig.12a), and in $\frac{B}{V}$, C , S (not shown here) have very similar tendencies as the functions of $\frac{m'}{m}$ in the above parameters (see Fig.10). However, the range of changes of the frequent processor number vs. V is narrower than the range of changes in $\frac{m'}{m}$. For example, in Fig.10a the number of used processors changes in range approx. $[0.4, 1]$. Here, the range

of changes is approx. $[0.3, 0.5]$ (Fig.12b). Even smaller ranges were observed in the experiments with changing $\frac{B}{V}, C, S$. It can be concluded that the size of frequent set of processors is growing with V , but not as quickly as the number of different used processors m' . It is because only a selected set of processors is frequently used while many other processors get to the solution due to the randomized selection.

In Fig.12c the load of the most frequently used processor is depicted vs. changing $A, \frac{B}{V}, C, S$. Independently of the type of changes the most loaded processor receives $0.6-0.75V$ on average. With growing A computation time starts dominating in schedule length, the processor selection method tends to build more computing power, and more processors are appended to the frequent set. Hence, the greatest piece of load sent to a single processor is diminishing. Growing $\frac{B}{V}$ allows for using fewer processors and for economizing on communication time. Hence, for big $\frac{B}{V}$ the most loaded processor receives almost $0.75V$. For small $\frac{B}{V}$ a big number of communications must be made anyway which expose the cost of communication startup times dominating in the schedule length. Consequently, GA becomes myopic to other processor parameters, the frequent set has more processors, and the load is more dispersed between the processors. The dependence on S is shown in Fig.12c is very weak. Note that this is an average over many sizes V . A more detailed picture exposes diversity with V similar to the one shown in Fig.11d, though in much narrower range. Unlike in Fig.11d the load sizes are generally decreasing with S , even for small loads V . Similarly to the results in Section 4.3 the biggest piece of the load received by a single processor does not depend on C .

The sum of the load assigned to all frequent processors is show in Fig.12d. As it can be seen the frequent processor set collects more than $0.8V$ on average. The function of the total load vs. $\frac{B}{V}$ has a minimum. This unexpected phenomenon can be explained in the following way. For big values of $\frac{B}{V}$ only a few processors take part in the computation because a single installment is sufficient to process the whole load. Therefore, the number of messages is small, load chunks have sizes close to processor memory buffer sizes, the frequent set has small cardinality and receives whole load. With decreasing $\frac{B}{V}$ more and more processors receive some load, and the contribution of the most loaded processor(s) is decreasing as depicted in Fig.12c. However, when $\frac{B}{V}$ becomes extremely small, communication cost is dominating schedule length, GA becomes unaware of processor parameters, and more of the processors are randomly included in the frequent set. Therefore, the cardinality of the frequent set is growing and also the total load in the frequent set is growing.

We finish the above exercise with a conclusion, that the frequent set of processors really exists. With the exception of the instances biased by small $\frac{B}{V}$ or big S , when almost all processors are frequent, the frequent set has approx. 40 – 50% of all available processors. They received 80-85% of the whole load, again with the exception of the cases biased by small $\frac{B}{V}$ or big S .

The results in Fig.11, and Fig.12 confirm that the set of processors receiving more load, and hence dominating in the computation exists. Yet, consider the

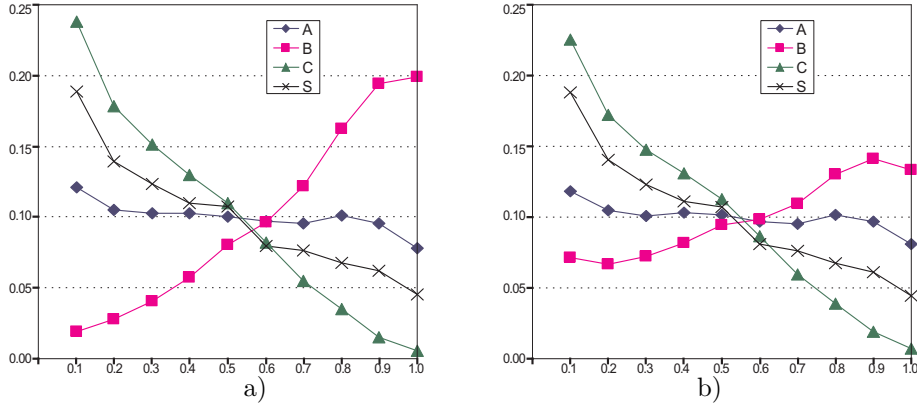


Figure 13: Load and number of messages vs processor rank in GA solutions. a) Load vs rank, b) Number of chunks vs rank.

method of test instance generation. When studying influence of a certain parameter, all processors have this parameter equal. We learned on the importance of the considered parameter via the consequences of its low, or high, values. But effects of the diversity of the given parameter were switched off. We did not verify how important this parameter may be if it had different values in the processor set. Therefore, a second set of 1000 instances were generated with $V = 100$, m generated from $U[1, 100]$, and A, B, C, S parameters generated from $U[0, 1]$. We examined the fraction of the whole load and the number of received messages against the rank of processors in the order of certain parameter value. The results of this study are shown in Fig.13.

In Fig.13 processors were grouped into sets comprising 10% of the processors ranked according to certain parameter. For example, value 0.3 on the horizontal axes in Fig.13 represent processors with relative rank $\frac{i}{m}$ in the range $(0.2, 0.3]$. The four functions depicted in Fig.13 correspond to four different rankings: according to A, B, C, S . Let us remind that for A, C, S smaller values represent better performance, and for B bigger values are better. The relationships are similar for the received load (Fig.13a) and for the number of messages (Fig.13b). Hence we will discuss only the load distribution. The distribution of the load is tightly connected with all processor parameters. It is evident that processors which have best communication links with respect to C , and S , the biggest memory buffers receive most of the load to process. The processors with small B , and big S, C , receive almost no load. For A the relationship is weaker but it is still noticeable (coefficient of correlation between A and the upper limit of rank box interval is ≈ -0.84).

We finish the study of the dominating set of processors with the following observations on the basis of computational experiments:

- the dominating processor set exists,

- the frequent processors set, as we defined it, comprises approximately 40-50% of all processors,
- in the biased case of big S , and small $\frac{B}{V}$ frequent processor set may include nearly all processors,
- there is a strong correlation between parameters C_i, S_i, B_i , and the amount of load received for processing.

4.5 Chunk Size Saturation

Another feature of problem solutions is load partitioning. After determining the sequence of communication and the overlaps, a linear program was used to find load distribution. Since the computational cost of linear programming may be considered high, it would be profitable to eliminate it in constructing good quality solutions. To examine the structure of load partitioning we analyzed the number of chunks which sizes equal the size of the target processor buffer, i.e. $\alpha_i = B_{\sigma(i)}$. We will call such chunks *full chunks*. It would be a very attractive solution to use just the processor buffer size as chunk size, thus eliminating the need for linear programming. The results of this exercise are shown in Fig.14.

In all figures shown in Fig.14 the number of full chunks is shown in relation to the total number of chunks n . The number of full chunks is almost always high or noticeable, but not all chunks are full. As it can be seen if Fig.14a,b,d, with growing size V the number of full chunks is also growing. This is intuitively reasonable because bigger load V requires more messages which expose costs of the startup time. These can be reduced by using as few messages as possible, and consequently filling the buffers more completely. This is also confirmed in Fig.14c where the number of full chunks is shown against changing $\frac{B}{V}$, and various values of V . When $\frac{B}{V}$ is small the number of messages must be big, hence the startup times dominate in the schedule length, and to reduce their contribution, the buffers are more fully filled. This situation is repeated in Fig.14d where the number of full chunks increases with the startup times. With growing A (Fig.14b) the number of full chunks is decreasing because the computation time starts dominating in the schedule length, not the startup times. Note that in Fig.14c the number of full chunks decreases with V , which may be attributed to the randomized nature of GA . With growing V greater number of messages must be sent. The message target processors are generated randomly. Hence with growing V chances are growing that a slow processor, or a processor with slow communication link may be selected to the communication sequence. The linear part tries to minimize influence of such bad choices by reducing chunk sizes. Hence with growing V the fraction of chunks which are not full also grows.

4.6 How Hard It Is To Find a Good Solution

In this section we study what makes our problem easy, or hard to solve. Let us introduce the goal of this section in more detail. Heuristics build good quality solutions for many combinatorial optimization problems. However, this good

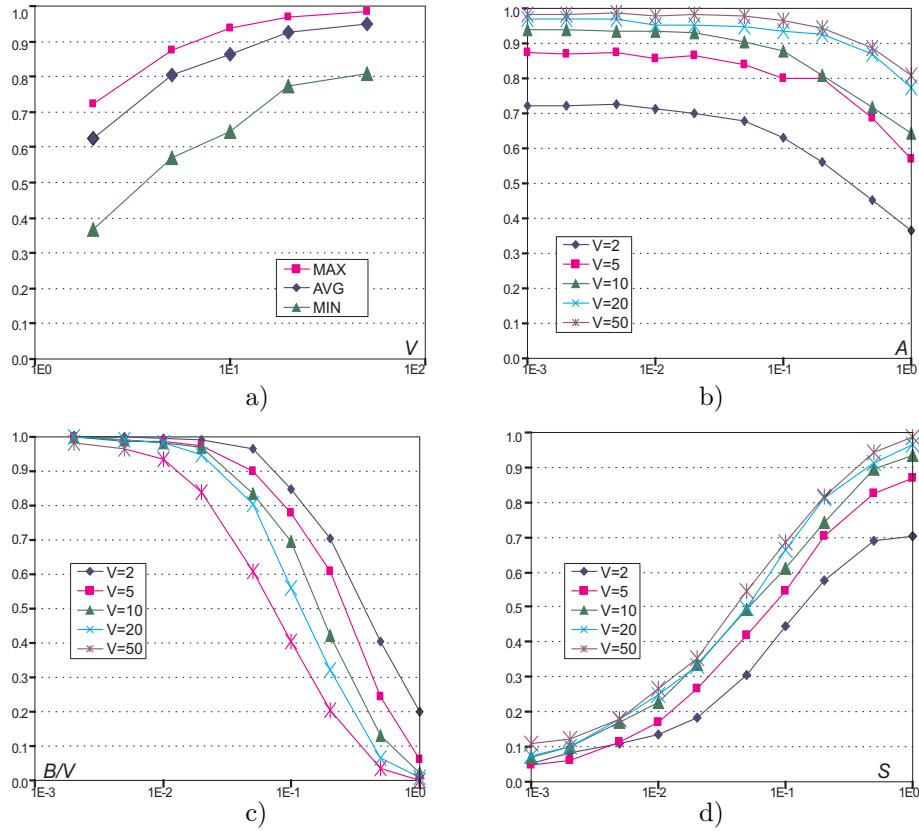


Figure 14: Average number of full chunks in GA solutions. a) vs V in experiments with changing A , b) vs A , c) vs $\frac{B}{V}$, d) vs S .

performance should be attributed to the nature of the problem, not a heuristic. Thus, it is possible that our genetic algorithm builds good solutions not because it is well designed, but because our scheduling problem may be easy to solve. If we learn which instances are easy, or hard, to solve then we will gain some new insights into the nature of the problem, and real merits of GA.

Now the problem is how to verify which instances are easy, and which ones are hard to solve. We will compare quality of the solutions obtained in three ways for various types of instances. The worst solution observed provides an indication on how bad a solution may be. The random solutions are not biased to being good or bad. GA solutions represent solutions which are optimized, and supposed to be good. The three solution types indicate what can be achieved in the worst case, without great efforts (random solutions), and at considerable cost of optimization. If GA solutions did not differ much from the random solutions, then it would signify bad GA design. All three algorithms were obtained using the GA infrastructure. The random solution is the best one in the initial GA

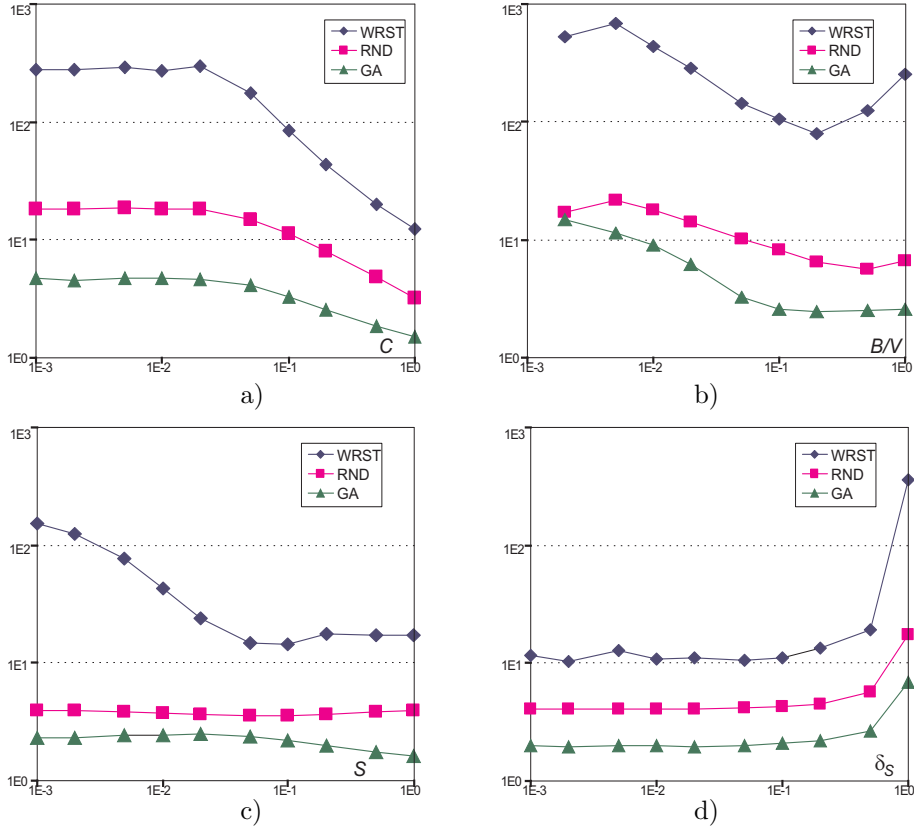


Figure 15: Quality of the solutions with reference to the lower bound for $V = 20$. a) vs C , b) vs $\frac{B}{V}$, c) vs S , d) vs the dispersion of S .

population of $G = 20$ solutions. The worst solution is the worst one observed in the course of solving some instance by GA. In all three algorithms linear programming was used to obtain the best chunk sizes α_i , and schedule length, for the given combinatorial part of the solution. Quality of the solutions is measured as the relative distance from the lower bound calculated in the following way. The minimum communication time is $\tau_1 = n_{MIN}S_{min} + VC_{min}$, where $C_{min} = \min_{i=1}^m C_i$, $S_{min} = \min_{i=1}^m C_i$. In this time at most $V_0 = (\tau_1 - S_{min}) \sum_{i=1}^m \frac{1}{A_i}$ load could be processed. The remaining load $V - V_0$ is processed in time at least equal to $\max\{0, \sum_{i=1}^m \frac{V}{1/A_i} - \tau_1 + S_{min}\}$. Thus the lower bound is equal to $\tau_1 + \max\{0, \sum_{i=1}^m \frac{V}{1/A_i} - \tau_1 + S_{min}\}$. We will examine performance of these three types of solutions for changing values and dispersion of system parameters.

In Fig.15 we have shown influence of the system parameters on the quality of the above three solution types. Fig.15a,b,c show results for the first set of random instances and $V = 20$. It is striking that the worst case solutions

(denoted *WRST*) can be over one order of magnitude further from the lower bound than the random solutions (denoted *RND*) or the solutions of the genetic algorithm (denoted *GA*). Moreover *GA* solutions are substantially better than *RND* solutions. Hence, *GA* really works. Now let us analyze the tendencies in Fig.15a,b,c. As it can be seen in Fig.15a with growing C all the lines tend to 1 which means that communication speed is decreasing, schedule length becomes dominated by the time of sending load off the originator. Hence in such a biased case it is getting easier to obtain good solutions. Similar tendency was observed for growing parameter A .

In Fig.15b the dependence of solution quality on changing $\frac{B}{V}$ is shown. With growing $\frac{B}{V}$ all the three types of solutions get closer to the lower bound. It is intuitively attractive to conclude that with growing $\frac{B}{V}$ good solutions are easier to obtain because we are less limited with the choice of the processor. Not disregarding this growing flexibility, it should not be forgotten that the construction of the lower bound influences the results presented here. The lower bound is based on the assumption that the smallest S_i coincide with the biggest B_i which is rarely true. Hence, for small $\frac{B}{V}$ and big number of the startups the error resulting from this simplification is significant. With increasing $\frac{B}{V}$ the domination of the startup costs in the schedule length decreases, and the lower bound is representing this situation better. Thus, the results in Fig.15b indeed confirm that with growing $\frac{B}{V}$ it is getting easier to obtain good quality solutions, however, it is achieved by using fewer messages and communication startup times. Moreover, for the biggest $\frac{B}{V}$ solutions *WRST*, *RND* are getting slightly worse and *GA* solutions are not. This means that even if memory buffers are big it is necessary to adjust the set of used processors. Genetic algorithm is doing it better than in the *RND* solutions.

In Fig.15c dependence of the three types of solutions on changing parameter S is shown. A counterintuitive tendency of improving *WRST* solution quality with growing S can be observed. With growing S the contribution of the startup time to schedule length is growing, independently of the chosen set of processors. Therefore, the difference between the worst solution and the lower bound is decreasing with growing S . Genetic algorithm is performing better than *RND* because it is able to build solutions with relative quality improving even with increasing domination of the startup time.

In Fig.15d quality of the solutions for growing dispersion of S is shown. The test instances for Fig.15d were generated as in the first set of instances with $V = 20$, except for parameter S which was generated with uniform distribution from range $[\frac{1-\delta_C}{2}, \frac{1+\delta_C}{2}]$. The value of δ_C is shown on the horizontal line in Fig.15d. As it can be seen with growing δ_C , and hence heterogeneity of the system, quality of all three types of solutions is worsening. This means that our problem is becoming harder to solve with growing heterogeneity of the computing environment. Similar experiments were performed for controlled dispersion $\delta_A, \delta_B, \delta_C$ of parameters $A, \frac{B}{V}, C$, respectively. In all cases the dependence of the quality of solutions on the range of diversity has very similar shape which once again confirms that in heterogeneous systems good quality solutions are

harder to obtain. Let us use the range of change of the worst-case solutions quality as an indicator of the sensitivity to the dispersion of certain parameter. For δ_S changing from 1E-3 to 1 the average distance from the lower bound grew ≈ 34 times. For similar changes of: 1) δ_C the distance changed ≈ 14 times, 2) δ_A changed ≈ 1.8 times, 3) δ_B changed ≈ 1.3 times. This means that diversity of S, C have the strongest influence on difficulty of obtaining good solutions, and diversity of $A, \frac{B}{V}$ the smallest.

We finish this section with the following conclusions.

1. It is easier to obtain good quality solutions when communication time or computation time dominate in the schedule length
2. It is easier to obtain good quality solutions for big memory buffers.
3. It is easier to obtain good solution quality for homogeneous systems. Solution quality is particularly sensitive to the dispersion of communication parameters S, C , and less to the dispersion of $A, \frac{B}{V}$.
4. Genetic algorithm really works, because it builds considerably better solutions than *RND*. Moreover, in some cases it is able to counteract the general tendencies of solution quality represented in *RND, WRST*.

5 Conclusions

In this paper we analyzed scheduling divisible loads on heterogeneous systems with communication startup times, limited memory buffers and multi- round load distribution.

In the first part of this paper a new, more realistic model of memory management was assumed. The problem of determining optimum communication schedule and load chunk sizes has been formulated as a mixed nonlinear program. This mixed nonlinear program has been reformulated for better solvability. Two parts in the solution can be distinguished: a combinatorial and an algebraic part. Two methods for solving the combinatorial part were proposed: branch and bound algorithm, and a genetic algorithm. The algebraic part is solvable by a linear program on condition that a solution from the combinatorial part is provided. The branch and bound algorithm turned out to be impracticable due to its prohibitive complexity. Therefore, in the following studies we relied on the solutions from the genetic algorithm.

In the second part of the paper we studied features of the solutions of our scheduling problem. The study was performed both analytically and by extensive computational experiments. Among the other, the following observations were made:

- It has been established that in the worst case arbitrarily big amount of load may have to be accumulated in the optimum solutions. By this feature multiprocessor schedules differ from the optimum schedules on a single processor. However, it turned out that in the near-optimum solutions obtained by the genetic algorithm accumulating the chunks is very rare.

- There is a minimum number of messages that must be sent to the processors anyway, it follows from the greatest processor buffer capacities. It can be shown that using this number of communications may result in arbitrarily bad solutions. In computational experiments it has been established that the number of messages is a small multiple of the minimum possible (1.4-10). Communication startup time is the main disincentive to using great number of messages in delivering the load to the processors.
- There is inequality in load distribution and a dominating set of processors receives most of the load. The size of the dominating set of processors is growing with size of the load V . There is a strong correlation between the parameters of a processor, and its contribution in load processing. Processors with faster communication links, bigger memory buffers, and computing faster receive more load. It appears that the order of parameter importance in load distribution is C_i, B_i, S_i, A_i .
- Majority of, though not all, load chunks carry maximum load, i.e. equal to the size of processor buffer. The number of full chunks grows with V , and is strongly correlated with S_i, B_i .
- The problem has natural tendency to become easier to solve when one parameter dominates in the schedule length. For example, big values of all A_i , and small C_i, S_i simplify obtaining good solutions.
- Another side of the above observation is that it is relatively easy to build biased instances which solutions are dictated by extreme values of certain parameter, e.g. extremely slow communication, or computation, or very small memory buffers.
- In a sense, parameters B_i and S_i go together when building a biased instance. Small memory buffers B_i incur many communications which expose cost of the startup time S_i . And vice versa, big startup times may be compensated by use of long messages which require big memory buffers.
- Good quality solutions are harder to obtain in heterogeneous systems.

We believe that the above set of observations may be helpful in constructing new, faster and still effective heuristics for the scheduling problem.

References

- [1] R. Agrawal, H.V. Jagadish, Partitioning Techniques for Large-Grained Parallelism, *IEEE Transactions on Computers* **37**(12), 1627-1634, 1988.
- [2] V.Bharadwaj, D.Ghose, V.Mani, T.Robertazzi, *Scheduling divisible loads in parallel and distributed systems*. IEEE Computer Society Press: Los Alamitos CA, 1996.

- [3] V.Bharadwaj, D.Ghose, V.Mani, Multi-installment Load Distribution in Tree Networks with Delays. *IEEE Transactions on Aerospace and Electronic Systems* **31**, 555-567, 1995.
- [4] Y.-C.Cheng, T.G.Robertazzi: Distributed computation with communication delay. *IEEE Transactions on Aerospace and Electronic Systems* **24**, 700-712, 1988.
- [5] M.Drozdowski, *Selected problems of scheduling tasks in multiprocessor computer systems*. Series: Monographs, No 321, Poznań University of Technology Press, 1997. <http://www.cs.put.poznan.pl/mdrozdowski/txt/h.ps>
- [6] M. Drozdowski, P. Wolniewicz, Divisible Load Scheduling in Systems with Limited Memory, *Cluster Computing* **6**(1), 19-29, 2003.
- [7] M.Drozdowski, M.Lawenda, Multi-installment Divisible Load Processing in Heterogeneous Systems with Limited Memory, in: R.Wyrzykowski *et al.*, PPAM 2005, *Lecture Notes in Computer Science* **3911**, 847-854, 2006.
- [8] M. Drozdowski, P. Wolniewicz, Optimum divisible load scheduling on heterogeneous stars with limited memory, *European Journal of Operational Research* **172**, 545-559, 2006.
- [9] X.Li, V.Bharadwaj, C.C.Ko, Processing divisible loads on single-level tree networks with buffer constraints, *IEEE Transactions on Aerospace and Electronic Systems* **36**, 1298-1308, 2000.
- [10] Lp_solve reference guide, 2007, <http://lpsolve.sourceforge.net/5.5/>.
- [11] T.Robertazzi, Ten reasons to use divisible load theory. *IEEE Computer* **36**, 63-68, 2003.
- [12] Y.Yang, H.Casanova, M.Drozdowski, M.Lawenda, and A.Legrand, On the complexity of multi-round divisible load scheduling, Research Report 6096, INRIA Rhône-Alpes, 38334 Montbonnot Saint Ismier, France, January 2007. <https://hal.inria.fr/inria-00123711>