

# Performance Modelling for Task-Parallel Programs

Matthias Kühnemann

Fakultät für Informatik  
Tech. Universität Chemnitz  
09107 Chemnitz  
Germany  
kuehnemann@informatik.tu-chemnitz.de

Thomas Rauber

Institut für Informatik  
Universität Halle-Wittenberg  
06099 Halle (Saale)  
Germany  
rauber@informatik.uni-halle.de

Gudula Rünger

Fakultät für Informatik  
Tech. Universität Chemnitz  
09107 Chemnitz  
Germany  
ruenger@informatik.tu-chemnitz.de

**Keywords:** execution time analysis, runtime formulas, mixed task and data parallelism, scientific computing.

## ABSTRACT

Many applications from scientific computing and physical simulations can benefit from a mixed task and data parallel implementation on parallel machines with a distributed memory organization, but it may also be the case that a pure data parallel implementation leads to faster execution times. Since the effort for writing a mixed task and data parallel implementation is large, it would be useful to have an a priori estimation of the possible benefits of such an implementation on a given parallel machine. In this article, we propose an estimation method for the execution time that is based on the modelling of computation and communication times by runtime formulas. The effect of concurrent message transmissions is captured by a contention factor for the specific target machine. To demonstrate the usefulness of the approach, we consider a complex method for the solution of ordinary differential equations with a potential for a mixed task and data parallel execution.

## 1. INTRODUCTION

We consider task parallel programs in a form where the entire program is built up from multi-processor tasks each of which can be executed on an arbitrary number of processors. Different multi-processor tasks can cooperate in a compositional way which means that one task produces data to be used as input data by another task. But a program may also contain independent tasks which can be executed concurrently to each other on subsets of processors, also called groups of processors. There usually exists a large variety of different parallel realizations of the same application program. The program versions can differ in the execution order of the tasks and the mapping of each multi-processor task onto a set of processors of a specific size. In this paper we address the question how to select an efficient parallel program version from the described class of programs when using a distributed memory machine (DMM).

Whether a pure data parallel or a mixed task and data parallel program version is more efficient strongly depends on the specific application program and its requirement for communication, especially for collective communication operations. Collective communication operations performed on smaller processor groups lead to smaller execution times due to the logarithmic or linear dependence of the communication times on the number of processors [?, ?, ?]. This behavior of the communication time can be an advantage for a concurrent computations on disjoint subsets

of processors. On the other hand, changing the data distribution from a pure data parallel execution to a group-oriented execution may require additional data redistribution which can outperform the advantage of a task parallel execution. In general, the effort to implement a task parallel program is considerable and so it would be useful to have an effective method to determine the most efficient program version before implementing all the details. Such a method needs a notion of costs assigned to different program versions under consideration where the costs of a specific version meet some requirements: they should be based on the specific task structure and they should reflect the behavior of the parallel execution times resulting on a specific parallel machine. A cost model with those properties is suitable for distinguishing between different program versions for which different task structures are known but final programs do not yet exist.

In this paper, we suggest a method based on costs estimating parallel execution times of task parallel programs with multi-processor tasks. The modelling of costs is based on a specification of the task structure and assumes a corresponding partition of the set of processors into groups or subsets of processors. The costs are given as runtime formulas that contain parameters describing the parallel target machine and characteristics of the algorithm to be implemented [?]. Runtime formulas have been used before for modelling the execution time of communication operations in isolation [?, ?, ?], but their use for modelling the execution time of task parallel programs has not yet been studied. The runtime formulas are built up from computation costs and communication costs which reflect the actual execution of a computation or a communication operation in isolation. But due to complex processor architectures, memory hierarchies and network properties, the execution time of an entire program might strongly be influenced by cache effects or network contention. Our aim is to investigate whether a purely program-oriented construction of the runtime formulas is suitable to model costs for task parallel executions and whether additional effects like cache effects and network contention have to be taken into consideration to obtain an accurate prediction.

In the following sections, we describe how to build up program-specific runtime formulas for parallel platforms with distributed memory. In Section 2 we summarize the general approach. Section 3 briefly introduces the parallel target platforms used, a Cray T3D and T3E and a Beowulf cluster of PCs, and presents the runtime formulas for communication operations for those machines. As complex example application we consider one-step methods for solving ordinary differential equations (ODEs) which have a potential for

multi-processor task parallelism. In Section 4, we model the runtime formula for an entire application. Section 5 shows that the runtime formulas reflect the performance behavior of the parallel program and Section 6 concludes the paper.

## 2. RUNTIME FORMULAS

The runtime formulas of our cost model are composed according to the structure of the programs whose execution time they describe. The runtime formulas consist of two parts describing the computation times and the communication times.

### 2.1 Computation operations

The computation times are modelled by taking into account the number of operations to be executed. For each arithmetic operation we use a computation time  $t_{op}$  that can be determined by runtime measurements with simple sequential test programs on the parallel machine used. In addition, we use execution times  $T_f$  for specific functions  $f$  which are needed to describe specific application problems. The execution time for such a function  $f$  can be build up from single computation operations or can be a measured time for the function in isolation. The reason to use a single value  $T_f$  in the runtime formula is that most numerical methods, like ODE solvers, are designed as black-box solvers, so that they can be used for arbitrary functions  $f$  describing the right-hand side of the ODE system. For different ODEs, the structure of the runtime formula remains the same but the specific cost  $T_f$  is exchanged when a different function  $f$  is used.

The memory access time may have a large influence on the execution time due to the temporal or spatial locality of the memory accesses of the specific program considered. Such effects can be captured by sophisticated approaches like the cache miss equations [?]. We could use this technique here to model the memory access time accurately and would get more complicated runtime formulas, but as it will be shown in Section 4, the simple approach of just using an average memory access time for the program does a sufficiently good job for regular programs from scientific computing. Moreover, when considering the modelling of communication times in the context of concurrent communication operations by different task parallel program parts, it is much more difficult to obtain accurate predictions since memory contention effects and the variations in the communication times dominate the variations in the execution times.

### 2.2 Communication operations

The communication times are modelled by formulas that describe the execution time of individual communication operations, such as single-transfer or broadcast operations. These formulas are used to model the internal communication time of a multi-processor task, which are often realized in a data parallel way. Communication formulas may also be used to describe the time for data redistributions inserted between different multi-processor tasks related in a compositional way. The runtime formulas for communication operations are functions in closed form that depend on the message size  $b$  and the number  $p$  of processors participating in the communication operation. For a given DMM, the value of these parameters result from modelling the runtime formulas by the least squares method. The execution of a broadcast operation, for example, on a Cray T3E using `MPI_Bcast()` can be described by the runtime formula

$$t_{sb}(p, b) = 97 \cdot p + 0.006 \cdot \log p \cdot b$$

In Section 3 we present runtime formulas for different communication operations in isolation, see also [?, ?, ?, ?].

### 2.3 Task parallel programs

Runtime formulas for entire programs are built up according to the task structure of the specific program. First, each participating multi-processor task  $M$  is assigned a runtime formula built up from communication and computation times according to the internal structure. We assume a data parallel or SPMD-like internal computation structure consisting of alternating phases of computation and communication, so that the execution time can be represented by

$$T_M(p, b) = T_{comp}(p, b) + T_{comm}(p, b).$$

$T_{comp}(p, b)$  is the maximum of the execution times of the participating processors.  $T_{comm}(p, b)$  is the sum of the runtimes of the communication operations used. Second, the runtime of an entire program is built up from the costs of the multi-processor tasks according to the task structure. In the case that two multi-processor tasks  $M_1$  and  $M_2$  are executed concurrently to each other on disjoint sets of processors, the maximum of the runtimes is taken:

$$\max(T_{M_1}(p_1, n_1), T_{M_2}(p_2, n_2)).$$

$T_{M_i}(p_i, n_i)$  denotes the runtime of task  $M_i$  for problem size  $n_i$  on a group of  $p_i$  processors,  $i=1,2$ . The sets of processors executing  $M_1$  and  $M_2$  are assumed to be disjoint. In the case that two multi-processor tasks  $M_1$  and  $M_2$  are executed one after another, the runtimes are added:

$$(T_{M_1}(p, n_1) + T_{M_2}(p, n_2)) + T_{redist}(p, m).$$

Both times  $T_{M_i}(p, n_i)$  contain the same parameter  $p$  for the sizes of the processor groups indicating that both tasks  $M_1$  and  $M_2$  are executed on a processor group of the same size, which is usually the same group. In addition, some communication between successive task activations might be necessary to realize redistribution of data. The time needed for the redistribution is denoted by  $T_{redist}(p, m)$  where  $p$  is the number of participating processors and  $m$  the size of the data to be redistributed.

### 2.4 Comparison with other cost models

Our approach of performance modelling with runtime formulas combines a modelling part for SPMD-like computations and a construction part for the upper level task parallelism of multi-processor tasks. The combination is inspired by the specific goal to compare the efficiency of different realizations with mixed task and data parallelism for the same given algorithm. Other cost models for parallel programming, like the the LogP model and its variations [?, ?] or the BSP model [?, ?, ?], are aimed at aspects different from ours. The LogP model is more architecture oriented with an emphasis on a detailed modelling of parallel runtimes but is less suited to express explicitly upper level task parallelism. The BSP model concentrates on a specific program structure of supersteps which is not straightforwardly related to a mixed task and data parallelism with multi-processor tasks. The advantage of using runtime formulas for the modelling of the execution times lies in the fact that they can easily be adapted to a mixed task and data parallel execution of program parts. Moreover, runtime formulas are able to capture the effect that the same communication operation may lead to different execution times, depending on whether other communication operations are simultaneously executed by concurrent processor groups or not. This

effect cannot easily be captured by other models like the BSP model [?, ?, ?] or the LogP model and its variations [?, ?] that assume a pure data-parallel execution of the program. In these models, it is difficult to capture runtime effects that occur only in a task parallel execution. Using runtime formulas, such effects can be captured, e.g., by contention factors introduced in Section 3.2.

In contrast, the presented approach of runtime formulas exactly matches this program structure and its realization with communication libraries like MPI or PVM on distributed memory machines by separating the two programming levels also in the cost model, i.e. by combining different methods of runtime estimation for the different kinds of parallelism. More precisely, the machine behavior and the realization of the communication and other runtime formulas are captured within some parameters of the modelling part as abstractly as possible while getting reasonable runtime estimations for the data parallel computations as it will be discussed in the next section. The further construction of more complex runtime formulas exactly mimics the task structure to be considered. Altogether, this results in a method that is ideally suited to be used by application programmers to compare the expected parallel runtime of different task structures without dealing with low-level architectural and implementation details.

### 3. COMMUNICATION COSTS

We consider the communication on different machine with distributed memory a Cray T3D, a Cray T3E and a Beowulf-Cluster. The T3D and T3E use a three-dimensional torus network. The six communication links of each node are able to simultaneously support hardware transfer rates of 600 MB/s for the T3E.

The Beowulf Cluster CLiC ('Chemnitzer Linux Cluster') is build up of 528 Pentium III processors clocked at 800 MHz. The processors are connected by two different networks, the communication network and the service network. Both are based on the fast-Ethernet-standard, i.e., the PEs can swap 100 MBit per second. The service network (Extreme Block Diamant) allows external access to the cluster. The communication network (Cisco Catalyst) is used for inter-process communication between the PEs.

We consider message-passing programs that are coded using the MPI standard [?]. On the CLiC, LAM MPI 6.3.2 is used.

#### 3.1 Communication operations in isolation

For single-transfer operations and collective communication operations we consider different variants and measure their runtime on the CLiC and on a Cray T3E-1200. For single-transfer operations, the runtimes are obtained on a 4-processor partition with processor 0 sending data to processor 3. The message sizes are between 2 KByte and 400 KByte. Except for some anomalies of the buffered single-to-single transfer on the Cray T3E, the runtimes increase linearly with the message size. On the Cray T3E-1200 the standard single-to-single transfer is the fastest operation. For the CLiC we have additionally implemented a piecewise single-to-single transfer operation that splits the message to be transmitted into pieces of 4KB and sends the pieces separately. This is the fastest single-transfer operation on the CLiC. For collective communication operations which may involve more than two processors, we consider single-broadcast operations (MPI\_Bcast()), accumulation operations (MPI\_Reduce()), gather operations (MPI\_Gather()), scatter operations (MPI\_Scatter()), and multi-broadcast oper-

ations (MPI\_Allgather()). The execution times for collective communication operations on the Cray are essentially faster than on the CLiC. Again, we have implemented piecewise communication operations on the CLiC that turn out to be faster than the original operations.

**Runtime formulas:** The runtime formulas that are used for the modelling are summarized in Table 1. The value  $b$  denotes the message size in bytes,  $p$  is the number of processors participating in the communication operation. The coefficients  $\tau$  and  $t_c$  can be considered as startup and byte-transfer times and are determined by curve fitting with the least-squares method. For different internal realization of the same communication operation, different values for the coefficients are obtained. We therefore use an additional parameter  $V$  to distinguish the different variants; the specific values for the coefficients  $\tau(V)$ ,  $t_c(V)$ ,  $t_1(V)$  and  $t_2(V)$  are given in Tables 2 - 3.

operation	runtime formula
MPI_Send	$t_{s2s}(b) = \tau + t_c \cdot b$
MPI_Bcast	$t_{sb}(p, b) = \tau \log_2(p) + t_c \cdot \log_2(p) \cdot b$
MPI_Reduce	$t_{sa}(p, b) = \tau \log_2(p) + t_c \cdot \log_2(p) \cdot b$
MPI_Allgather	$t_{mb}(p, b) = \tau_1 + \tau_2 \cdot p + t_c \cdot p \cdot b$

**Table 1:** Runtime formulas for communication operations.

**Single transfer:** The runtime for a single-to-single transfer is modelled by a linear function  $t_{s2s}(b) = \tau + t_c \cdot b$  where  $\tau(V)$  denotes a *startup time* and  $t_c(V)$  denotes the *byte-transfer time* for the specific operation  $V$ . Curve fitting results in the parameter values given in Table 2. The negative coefficients arise when using the least squares method for measured runtimes in the complete range of message sizes. Restricting the method to small message sizes up to 2 KByte leads to a separate, more accurate runtime formula for small messages, which does not contain negative coefficients. A comparison between the predicted and measured runtimes for all MPI single-transfer operations shows that the predictions fit the measured runtimes quite accurately, especially for large messages (not shown in a figure).

variant	Coefficient for runtime formula of single-transfer, single-broadcast and accumulation			
	CLiC		Cray T3E	
	$\tau(V)[\mu s]$	$t_c(V)[\mu s]$	$\tau(V)[\mu s]$	$t_c(V)[\mu s]$
MPI_Send	-7.526	0.10607	13.965	0.00267
Send_P	145.021	0.08804	9.748	0.00533
MPI_Bcast	-3420.688	0.3409	7.723	0.0039
Bcast_P	564.125	0.0939	7.007	0.0050
MPI_Reduce	-119.871	0.1223	168.516	0.0093

**Table 2:** Parameter values for single-transfer operations and for the runtime of collective operations with a logarithmic dependence on the number of participating processors.

**Single broadcast and accumulation:** The modelling of single-broadcast (MPI\_Bcast()) operations is performed with the formula from Table 1. The logarithmic dependence on the number  $p$  of processors is used because the broadcast transmissions are based on broadcast trees with logarithmic depth. The same formula can also be used for the prediction of single-accumulation operations (MPI\_Reduce()). Again, the values of the parameters  $\tau(V)$  and  $t_c(V)$  are determined from a set of measured runtimes.

**Multi-broadcast and gather/scatter:** The predictions for multi-broadcast operations (`MPI_Allgather()`) are made using the formula from Table 1, which increases linearly with the message size and the number of processors. This formula can also be used for the prediction of scatter and gather operations. The resulting coefficients are shown in Table 3.

Coefficient for multi-broadcast/gather/scatter						
variant	CLiC [ $\mu s$ ]			Cray T3E [ $\mu s$ ]		
	$\tau_1(V)$	$\tau_2(V)$	$t_c(V)$	$\tau_1(V)$	$\tau_2(V)$	$t_c(V)$
MultiBcast	-33270.8	21801.5	1.031	-3.72	42.60	0.028
MultiBcast_P	-9724.2	8385.5	0.690	-24.20	-0.81	0.036
MPI_Allgather	9175.3	-7542.0	3.182	6.04	-0.75	0.019
Allgather_P	-669.9	543.3	2.806	7.72	13.23	0.018
MPI_AllgatherV	-709.6	-957.0	5.443	11.83	17.47	0.019
MPI_Gather	-316.2	654.5	0.095	31.08	-118.51	0.006
MPI_Scatter	24.6	1439.5	0.086	-0.48	5.45	0.003

**Table 3:** Parameter values for the runtime formulas with a linear dependence on the number of participating processors.

### 3.2 Communication in task parallel executions

In the last subsection, we have considered runtime formulas for the execution time of communication operations in isolation, which means that there was no concurrent communication operation of the same application program in transmission. In Section 4 those runtime formulas are used for the modelling of data parallel realizations of regular programs from scientific computing [?, ?]. But it is not a priori clear whether those runtime formulas (where as the parameter  $p$  the group size performing a communication is taken) can also be used for modelling the execution time of mixed task and data parallel programs and in fact, experiments have shown that using the runtime formulas from Subsection 3.1 leads to predicted runtimes that are too small compared to the measured runtimes. The effect is much larger on the Beowulf cluster than on the T3E. This behavior indicates that concurrent communication operations of the same application program can interfere with each other in the sense that the transmission time is delayed, especially when many single-to-single transmissions are executed concurrently.

Further experiments have shown that only a slight change of the runtime formula for communication operations is needed to model the runtime of concurrent communication transmissions. The change can be concluded from the following observation made when using the runtime formulas for data parallel realizations from Subsection 3.1 to model the runtime for concurrent collective communication operations: The difference between the measured and the predicted runtimes is increasing with the number  $p$  of processors (used for the entire application program) and the size  $n$  of the message transmitted. The dependence on  $p$  is more significant than the dependence on  $n$ . Both dependencies correspond to the expectation that collisions in the network become more likely when the number of participating processors and the message sizes are increasing. To capture the interference of a concurrent execution of communication operations in the runtime formulas, we adjust the byte transfer time of the runtime formulas by introducing a network contention factor  $C(p, n)$  that depends on  $p$  and  $n$ . In principle, this factor can be determined by a detailed queuing analysis, but this would require the knowledge of the internal realizations of the MPI communication operations. Thus, we take the approach to determine the contention factor empirically starting from the runtime estimation with the runtime formulas from Subsection 3.1.

We determine the contention factor by comparing the delay in the execution times for communication operations, when they are performed concurrently, with the execution times of these operations without network contention. The contention factor itself is modelled as a function of  $p$  and  $n$  and the shape of the contention factor (or contention function) has been determined by experiments from the measured delays in the execution times. The coefficients within this function are then determined by curve fitting. The resulting contention factor is used for the modelling of the communication times of those program parts of complex application programs that are executed in a mixed task and data parallel way. For each parallel machine, a different contention factor results because of the different network architectures of the machines. To summarize, using the contention factor, the structure of the runtime formulas for one communication operation does not change but the contention factor is used to adjust the byte-transfer time  $t_c$  to the specific situation. Especially, the startup times are not changed. Thus, for example, the runtime formula of a multi-broadcast operation is

$$\tilde{t}_{mb}(p, b) = \tau_1 + \tau_2 \cdot p + C(p, n) \cdot t_c \cdot p \cdot b. \quad (1)$$

The same contention factor is used for each communication operation, i.e., there is no need to determine different contention factors for different communication operations. The communication operations of pure data parallel parts of application programs are modelled without the contention factor because there is no interference of message transmission. The contention factors for the modelling of concurrent message transmissions are

$$\begin{aligned} C_{T3D}(p, n) &= 0.07 \cdot p \cdot \log_2(\log_2(\log_2(p))) \cdot \log_2(n) \\ C_{T3E}(p, n) &= 0.04 \cdot p \cdot \log_2(\log_2(p)) \cdot \log_2(n) \\ C_{CLiC}(p, n) &= 0.0045 \cdot p \cdot p \cdot \log_2(p) \cdot (\log_2(n) + \log_2(p)), \end{aligned}$$

where  $p$  denotes the total number of processors participating in the execution and  $n$  the size of the message transmitted.

## 4. COMPLEX EXAMPLE APPLICATION

To investigate the usefulness of the modelling approach for complex application programs, we consider parallel implementations of a specific solution method for ordinary differential equations (ODEs), the iterated Runge-Kutta method (iterated RK method).

### 4.1 Task structure of iterated RK methods

The iterated RK method is an explicit one-step method for the solution of initial value problems of ODEs. The iterated RK methods determines a sequence of approximation values  $y_1, y_2, y_3 \dots$  for the exact solution of the ODE System in a series of sequential time steps. In each of the time steps a fixed number  $s$  of stage vectors are iteratively computed and combined to the next approximation vector in the following way:

for  $l = 1, \dots, s$  initialize stage vector  $v^l(0)$   
for  $j = 1, \dots, m$   
for  $l = 1, \dots, s$ : compute new stage vector approximation  $v^l(j)$   
compute new approximation vector  $y_{k+1}^{(m)}$ ,  
compute approximation vector  $y_{k+1}^{(m-1)}$  for step size control.

The number  $m$  of iterations is given by the specific RK method.

Each computation of a stage vector approximation requires an evaluation of the function  $f$  that describes the ODE to be solved. The advantage of the iterated RK methods for a parallel execution is that the iteration system of size  $s \cdot n$  consists of  $s$  independent function evaluations that can be performed in parallel [?]. For systems of differential equations, an additional data parallelism can be exploited. Thus, the algorithm provides several possibilities for a parallel implementation. The computation of the stage vectors in one iteration  $j$  of the stage-vector computation can be performed on subsets of processors at the same time (*group implementation*) or alternatively by all processors one after another (*consecutive implementation*). The group implementation is a mixed task and data parallel implementation whereas the consecutive implementation is a pure data parallel realization.

## 4.2 Consecutive implementation

The computation time of the consecutive execution on  $p$  processors without the stepsize control can be modelled by the formula

$$T_{dp}(n, p) = \left( ms \left\lceil \frac{n}{p} \right\rceil + \left\lceil \frac{n}{p} \right\rceil \right) (2s+1) t_{op} + \left( ms \left\lceil \frac{n}{p} \right\rceil + s \left\lceil \frac{n}{p} \right\rceil \right) T_f + n s t_{op}$$

where  $T_f$  denotes the time for the evaluation of *one* component of  $f$  and  $t_{op}$  denotes the time for the execution of one arithmetic operation. In each loop body each processor has to compute  $n/p$  components of the argument vector using  $2s + 1$  operations and  $n/p$  components of  $f$ . Since  $f$  and its access pattern are not known in advance, the complete argument vector has to be made available to each processor with a multi-broadcast operation. For the communication time, the following formula is used with  $t_{mb}$  from Table 1.

$$C_{dp}(n, p) = s m t_{mb} \left( p, \left\lceil \frac{n}{p} \right\rceil \right) + t_{mb} \left( p, \left\lceil \frac{n}{p} \right\rceil \right).$$

## 4.3 Group implementation

A group implementation of the iterated RK method uses  $s$  independent groups of processors where each group  $G_i$  of size  $g_i$ ,  $i = 1, \dots, s$ , is responsible for computing the approximations of one specific stage vector. The pseudocode is illustrated in Figure 1. The processor groups should be of equal size, since the computation of each stage vector approximation requires an equal amount of computation. As it is possible that  $p$  is not a multiple of  $s$ , the group with the smallest number  $g_{min} = \lfloor p/s \rfloor$  of processors determines the computation time

$$T_{gp}(n, p) = \left( m \left\lceil \frac{n}{g_{min}} \right\rceil + \left\lceil \frac{n}{p} \right\rceil \right) (2s+1) t_{op} + \left( m \left\lceil \frac{n}{g_{min}} \right\rceil + \left\lceil \frac{n}{g_{min}} \right\rceil \right) T_f + \left\lceil \frac{n}{g_{min}} \right\rceil s t_{op}.$$

The communication time is modelled by the runtime formula with the machine specific contention factor, see Equation (1), to reflect the concurrently executed multi-broadcast operations.

$$C_{gp}(n, p) = 2 \cdot m \cdot \tilde{t}_{mb} \left( g_{min}, \left\lceil \frac{n}{g_{min}} \right\rceil \right) + t_{mb} \left( p, \left\lceil \frac{n}{p} \right\rceil \right).$$

```

forall  $l \in \{1, \dots, s\}$  do in parallel (group parallelism)
  forall processors  $q \in G_l$  do (data parallelism) {
    compute  $\lceil n/g_l \rceil$  components of  $f(y_k)$ ;
    initialize  $\lceil n/g_l \rceil$  components of  $\mu_{(0)}^1, \dots, \mu_{(0)}^s$ ;
  }
for  $j = 1, \dots, m$  do (sequential iteration) {
  forall  $l \in \{1, \dots, s\}$  do in parallel (group parallelism)
    forall  $q \in G_l$  do (data parallelism inside groups) {
      compute  $\lceil n/g_l \rceil$  components of argument  $\mu(l, j)$ ;
      group-multi-broadcast of local components of  $\mu(l, j)$ ;
      evaluate  $\lceil n/g_l \rceil$  components of  $f(\mu(l, j))$ ;
      group-multi-broadcast of local components of  $f(\mu(l, j))$ ;
    }
  forall processors  $q$  do (data parallelism on all processors)
    compute  $\lceil n/p \rceil$  components of  $y_{k+1}$ ;
    multi-broadcast the computed components of  $y_{k+1}$ ;
}
stepsize control;

```

**Figure 1:** Group implementation for one time step of the iterated RK method.

## 5. RUNTIME EXPERIMENTS

To validate the runtime formulas we have performed runtime tests on the T3D and T3E and the CLiC for up to 128 processors. Since the execution time of the iterated RK method strongly depends on the ODE system to be solved, we consider two classes of ODE systems:

*Sparse ODE systems.* These ODE systems have a right-hand side function  $f$  for which the evaluation of each component has a fixed evaluation time that is independent of the size  $n$  of the ODE system (sparse function), i.e., the evaluation time of the entire function  $f$  consisting of  $n$  components increases linearly with the size of the ODE system.

*Dense ODE systems.* These ODE systems have a right-hand side function  $f$  for which the evaluation of each component has an evaluation time that increases linearly with  $n$ , i.e., the evaluation time of the entire function  $f$  increases quadratically with the size of the ODE system.

Dense ODE systems lead to the fact that the computation time usually dominates the communication time, i.e., the comparison of measured and predicted execution times shows how good the computation times are modelled. On the other hand, for sparse ODE systems the communication time plays an important role and the comparison also includes this part. We have used an iterated RK method with  $s = 4$  stages that is based on an implicit RadauIIA method. This method leads to a convergence order of 7, if  $m = 6$  iterations are executed in each time step. The runtimes shown in the following tables and figures are runtimes for one time step of the method and are obtained by averaging over a large number of time steps. The (measured and predicted) runtimes include the time for stepsize and error control.

On the T3D and T3E, the runtime of the *consecutive* implementations of the iterated RK method for *dense* ODE systems can be modelled very accurately and are not shown here. But sparse ODE systems also lead to good predictions, see Table 4. In this case, no concurrent message transmissions take place, and therefore the contention factor does not need to be used. Table 5 evaluates the

$n$	measured runtimes				predicted runtimes			
	p=32	p=64	p=128	p=256	p=32	p=64	p=128	p=256
242	0.017	0.028	0.053	0.104	0.020	0.037	0.075	0.163
882	0.026	0.034	0.058	0.108	0.028	0.043	0.080	0.168
1922	0.041	0.045	0.067	0.115	0.042	0.053	0.067	0.174
3362	0.062	0.061	0.078	0.127	0.060	0.066	0.089	0.184
5202	0.089	0.082	0.094	0.140	0.084	0.083	0.099	0.196
7442	0.123	0.107	0.114	0.157	0.113	0.103	0.112	0.211
10082	0.163	0.137	0.139	0.179	0.147	0.128	0.148	0.227
13122	0.206	0.172	0.170	0.206	0.187	0.156	0.170	0.246
16562	0.257	0.211	0.203	0.239	0.230	0.187	0.196	0.268

**Table 4:** Consecutive implementation for *sparse* ODE systems, on a T3D. The average deviation is 10.3 %.

$n$	measured runtimes				predicted runtimes			
	p=32	p=64	p=128	p=256	p=32	p=64	p=128	p=256
242	0.013	0.019	0.039	0.096	0.023	0.040	0.077	0.156
882	0.036	0.056	0.106	0.218	0.036	0.057	0.106	0.216
1922	0.063	0.088	0.164	0.308	0.058	0.086	0.158	0.319
3362	0.103	0.138	0.253	0.450	0.090	0.129	0.236	0.474
5202	0.153	0.203	0.363	0.642	0.131	0.187	0.338	0.679
7442	0.212	0.281	0.494	0.886	0.183	0.259	0.470	0.937
10082	0.286	0.376	0.665	1.166	0.245	0.346	0.628	1.256
13122	0.368	0.484	0.854	1.489	0.317	0.449	0.816	1.629
16562	0.466	0.607	1.068	1.870	0.400	0.566	1.032	2.061

**Table 5:** Group implementation for *sparse* ODE system on a T3D. The average deviation is 13.5 %.

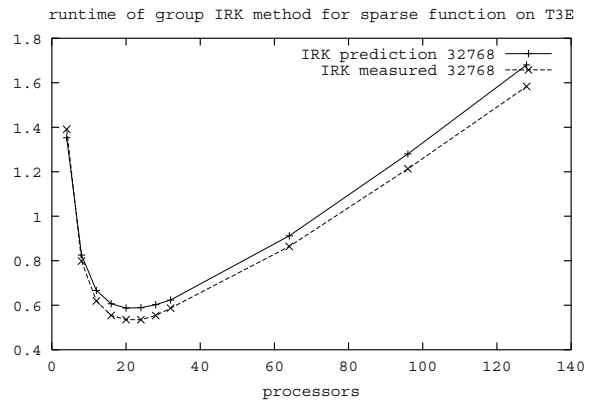
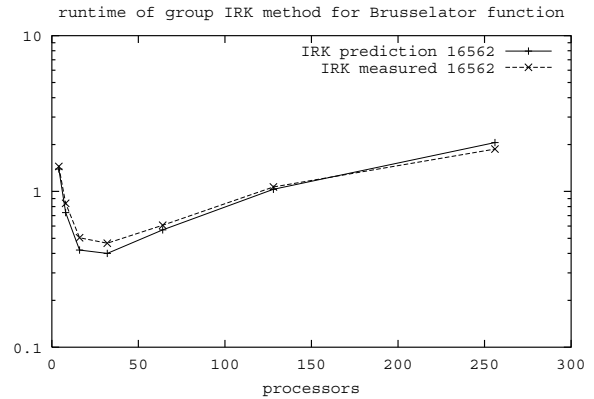
prediction for the *group* implementation where the contention factor is used. For both cases, the predictions are quite accurate, but not as accurate as the predictions for dense ODE systems. Nevertheless, they are accurate enough to be used for predicting the effect of a task parallel implementation. Table 6 shows the predictions for a *group* implementation of the iterated RK method for *dense* ODE systems on a T3E, again using the contention factor for the concurrent message transmissions because of a task parallel execution. Table 7 shows the same predictions for *sparse* ODE systems. Figure 2 illustrates the accuracy of the *group* implementation for *sparse* ODE systems of large sizes on the T3D and T3E for different numbers of processors. In contrast to dense ODE systems, the solution of sparse ODE systems leads only to considerable speedups for processor numbers of up to 16. For larger numbers of processors, the communication time dominates the computation times.

$n$	measured runtimes				predicted runtimes			
	p=16	p=32	p=64	p=128	p=16	p=32	p=64	p=128
100	0.012	0.012	0.017	0.035	0.016	0.015	0.020	0.032
500	0.194	0.108	0.074	0.082	0.210	0.116	0.076	0.069
1000	0.753	0.394	0.227	0.176	0.808	0.419	0.235	0.164
1500	1.682	0.862	0.468	0.310	1.801	0.920	0.492	0.305
2000	3.227	1.639	0.865	0.524	3.187	1.616	0.848	0.502
3000	7.041	3.555	1.840	1.043	7.145	3.602	1.860	1.038
5000	19.905	9.992	5.097	2.746	19.789	9.946	5.070	2.718

**Table 6:** Group implementation for *dense* ODE system on the T3E. The average deviation is 6.4 %.

$n$	measured runtimes				predicted runtimes			
	p=16	p=32	p=64	p=128	p=16	p=32	p=64	p=128
32	0.004	0.006	0.014	0.031	0.007	0.010	0.017	0.028
128	0.005	0.008	0.016	0.036	0.009	0.012	0.018	0.031
288	0.008	0.011	0.020	0.043	0.011	0.014	0.021	0.036
800	0.018	0.022	0.037	0.075	0.020	0.022	0.032	0.055
1568	0.030	0.035	0.055	0.106	0.032	0.035	0.049	0.086
3200	0.058	0.064	0.098	0.186	0.061	0.063	0.089	0.156
6272	0.109	0.118	0.178	0.329	0.116	0.117	0.167	0.299
15488	0.261	0.280	0.416	0.763	0.284	0.289	0.416	0.758
32768	0.554	0.587	0.864	1.583	0.607	0.624	0.912	1.680

**Table 7:** Group implementation for *sparse* ODE system on the T3E. The average deviation is 16.5 %.



**Figure 2:** Measured and predicted execution times for a group implementation of one time step of the iterated RK method for a Brusselator system of size  $n = 16562$  on a T3D (top) and of size  $n = 32768$  on a T3E (bottom).

Table 8 compares measured and predicted execution times for *group* implementations of the iterated RK methods for *dense* ODE systems on the CLiC. Figure 3 shows an illustration for a fixed message size and different numbers of processors. For this machine, the deviations between the measured and predicted execution times are much larger than for the T3D and T3E, especially for a larger number of processors. Although the deviations may be quite large, they can still be used as a rough estimate of the performance of a task parallel execution. The main reason for the large deviations are caused

$n$	measured runtimes				predicted runtimes			
	p=16	p=32	p=64	p=128	p=16	p=32	p=64	p=128
100	0.018	0.030	0.060	0.494	0.011	0.022	0.082	0.424
500	0.104	0.129	0.216	2.423	0.092	0.127	0.406	1.797
1000	0.301	0.296	1.288	3.425	0.283	0.306	0.840	3.667
1500	0.600	0.515	2.274	3.446	0.572	0.535	1.302	5.451
2000	1.049	1.492	2.418	4.296	0.959	0.813	1.790	7.372
3000	2.313	2.656	2.506	3.601	2.033	1.519	2.856	11.154
5000	7.814	5.180	4.714	15.303	5.361	3.528	5.289	19.036

**Table 8:** Group implementation for *dense* ODE system on the CLiC. The average deviation is 25.5 %.

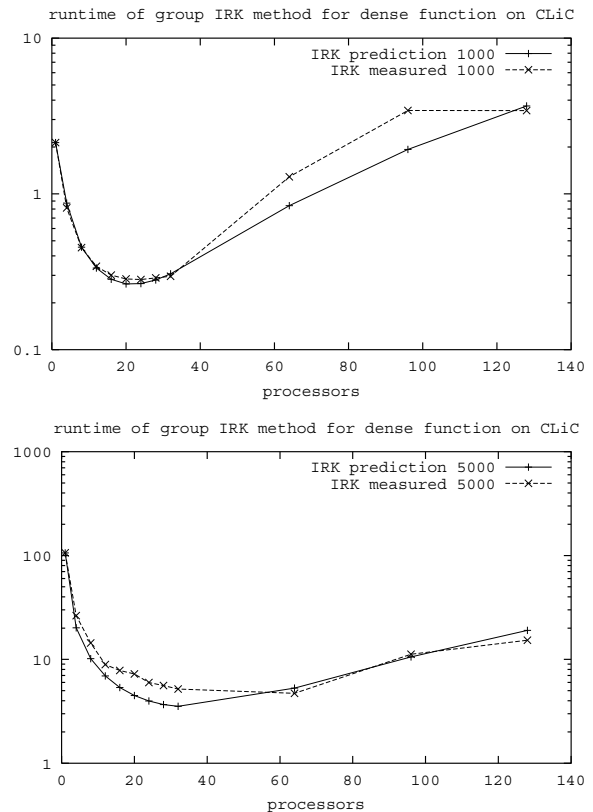
by the large increase of the communication time with an increasing number of processors. But also for a smaller number of processors, there are considerable deviations for large system sizes because of caching effects which are caused by the memory hierarchy of the single processors (Intel Pentium III). Such effects are much smaller on the Alpha 21164 processor of the T3E-1200 because of its different cache organization, see [?] for a more detailed investigation of such effects. The runtimes on the CLiC show that even for dense ODE systems, the machine only leads to satisfactory speedups for up to 32 processors. For larger processor numbers, the communication time and the network contention are too high.

## 6. CONCLUSIONS

In this article, we have shown that it is possible to model the execution times of mixed task and data parallel implementations by runtime formulas and that the use of a simple contention factor is sufficient to capture the interference of concurrent message transmissions. The runtime formulas model the execution times quite accurately for parallel machines like the T3D and T3E with a high-speed interconnection network. For a Beowulf cluster with an Ethernet-based network, the network contention caused by concurrent transmissions is much larger and it is more difficult to capture the effects by a simple contention factor. But the predictions are still reasonable and give a first impression of the possible effects of a task parallel realization.

## Acknowledgment

We thank the NIC Jülich for providing access to a Cray T3E and the EPCC Edinburgh for providing access to a Cray T3D.



**Figure 3:** Measured and predicted execution times for a group implementation of one time step of the iterated RK method for dense ODE of size  $n = 1000$  (top) and  $n = 5000$  (bottom) on the CLiC Beowulf cluster.