

## A FUNCTIONAL LANGUAGE FOR DEPARTMENTAL METACOMPUTING

FRÉDÉRIC GAVA and FRÉDÉRIC LOULERGUE  
*Laboratory of Algorithms, Complexity and Logic  
Créteil, France; {gava,loulergue}@univ-paris12.fr*

Received 16 September 2004

Revised 31 October 2004

Communicated by Sergei Gorlatch

### ABSTRACT

We have designed a functional data-parallel language called BSML for programming bulk synchronous parallel (BSP) algorithms. Deadlocks and indeterminism are avoided and the execution time can be then estimated. For very large scale applications more than one parallel machine could be needed. One speaks about metacomputing. A major problem in programming application for such architectures is their hierarchical network structures: latency and bandwidth of the network between parallel nodes could be orders of magnitude worse than those inside a parallel node. Here we consider how to extend both the BSP model and BSML, well-suited for parallel computing, in order to obtain a model and a functional language suitable for metacomputing.

*Keywords:* Functional Programming, Bulk Synchronous Parallelism, Metacomputing.

### 1. Introduction

Some problems require performance that can only be provided by massively parallel computers. For very large scale applications more than one parallel machine could be needed. One speaks about *metacomputing* [18]. In recent years there has been a trend towards using a set of parallel machine for these kinds of problems. Metacomputing infrastructures couple multiple clusters or parallel machines via a wide-area network. Research on global computational infrastructures has raised considerable interest in running parallel applications on *distributed* systems. Programming this kind of *metacomputers* is still difficult due to the presence of different networks, local and wide-area ones. High-level language and *cost models* are needed to ease the programming of *hierarchical* architectures such as *clusters of clusters*.

Bulk-Synchronous Parallel ML or *BSML* is an extension of ML for programming *direct-mode* Bulk Synchronous Parallel (BSP) algorithms as functional programs. BSP computing is a parallel programming model introduced by Valiant [19] to offer a high degree of abstraction like PRAM models and, yet, allow *portable* and *predictable* performance on a wide variety of architectures. BSML expresses them with a small set of primitives taken from the *confluent*  $BS\lambda$ -calculus [13].

But metacomputing programs need a more detailed model, including latency and bandwidth of the local and wide-area (or intranet) networks, the number of

clusters or parallel machines and the number of processors in each parallel machine. We currently make some simplifying assumptions about the networks: we use stable topologies, latencies and bandwidth. When the parallel machines are still in the same organization, university or building, this kind of programming is usually called *departmental metacomputing* [1]. In this way, regular network performances are certainly realistic for the duration of a part of a program and are still less sensitive to security measures. Assuming a metacomputer with stable network performances allows us to focus on the impact of network performance to design a cost model for this kind of architecture and a functional language for a *high-level* programming point of view. Our ultimate goal is to develop a functional language which could go beyond these limitations.

In [8] a previous version of the DMM cost model and the primitives of DMML were presented, but without a parallel implementation, a formal semantics, examples and experiments. A formal semantics is the most precise specification one could give for a new programming language. In our case it was a necessary very valuable guide for the parallel implementation. We also proved that the DMML language is deterministic.

Such a semantics is also a requirement for proving the correctness of programs. As usual in functional languages, we could *prove* the correctness of the DMML implementations of DMM algorithms with a *proof assistant* as done in [7] for BSP algorithms implemented in BSML. Using the *extraction* capability of proof assistants, we could generate a *certified* implementation of these algorithms. Nevertheless proofs of correctness of parallel programs (and also of the compilers and runtimes systems) go far beyond the scope of this paper and constitute the Propac project ([www.propac.free.fr](http://www.propac.free.fr)).

In section 2 we briefly review the BSP model and how to extend it for departmental metacomputing by adding a new level of communication. Then, we present informally our new functional parallel language, called Departmental Metacomputing ML or DMML (Section 3), a formal semantics (Section 4) for a core sub-language. Section 5 is devoted to two examples of collective communication operations and to the implementation. We discuss related work (Section 6) and conclude (Section 7).

## 2. A Model for Departmental Metacomputing

We assume throughout this paper that a *metacomputer* is a set of multiple clusters or parallel machines, i.e. a *cluster of clusters*, with fully connected local networks (LAN) and a fully connected intranet network, here excessively called WAN or *departmental* WAN. Each parallel machine has a gateway that connects its private LAN to the WAN. Here we give the name of *departmental metacomputing* system because the parallel machines are within the same organization. In this way, we use regular topologies, constant latencies and bandwidths.

### 2.1. Bulk Synchronous Parallelism

A BSP computer contains a set of uniform *processor-memory* pairs, a *communication network* allowing inter-processor delivery of messages and a *global synchronization unit* which executes collective requests for a *synchronization barrier*. For the sake of

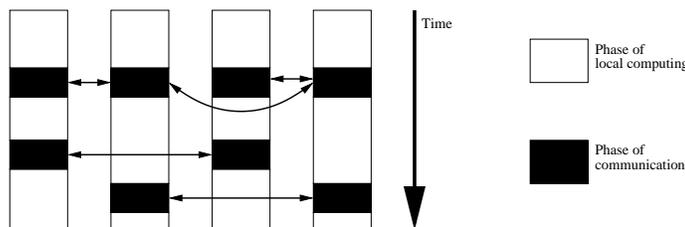


Figure 1: The MPM model of computation

conciseness, we refer to [17] for more details. In this model, a parallel computation is divided in *super-steps*, at the end of which exchanges of data and a synchronization barrier are performed. Hereafter all requests for data which have been posted during a preceding super-step are fulfilled.

The performance of the machine is characterized by 3 parameters:  $p$  is the number of processor-memory pairs,  $l$  is the time required for a global synchronization and  $g$  is the time for collectively delivering a 1-relation (communication phase where every processor receives/sends at most one word).  $g$  and  $l$  are expressed as multiples of the local processing speed  $s$ . The network can deliver an  $h$ -relation in time  $g \times h$  for any arity  $h$ . The execution time of a super-step is thus the sum of the maximal local processing time, of the data delivery time and of the synchronization time.

## 2.2. Discussion about the BSP model

There are two main arguments against using BSP for metacomputing. First the global synchronization barrier is claimed to be expensive especially for a set of parallel machines. Second, this model does not take into account the different capacities of the parallel machines and different networks: it is not a heterogeneous and hierarchical model of computation. Our proposal attempts to give a solution to the aforementioned problems. Starting from BSP, we address the problem of enlarging the number of parameters without introducing an unbearable complexity.

To remedy the first problem, [3] introduces the MPM model of computation which is a model directly inspired by the BSP model. It offers to replace the notion of super-step by the notion of *m-step* defined as: at each m-step, each process performs a sequential computation phase, then a communication phase (Figure 1), black boxes are phases of communication. During these communication phases, the processes exchange the data they need for the next m-step.

To remedy the second problem, [14] investigated a two-level hierarchical BSP model for a cluster of SMP machines and two-levels of communications. A BSP<sup>2</sup> computer consists of a number of *uniformly* BSP units, connected by a communication network. Execution of a BSP<sup>2</sup> program proceeds in *hyper-steps* separated by global synchronizations. On each hyper-step each BSP unit performs a complete BSP computation and then communicates with other BSP units. However, the authors noted that none of the algorithms they have analyzed shows any significant benefit from this approach and the experiments do not follow the model. The failure of the BSP<sup>2</sup> model to provide any major performance comes from three main reasons: first the BSP units are generally different in practice, second the synchro-



```

bsp_p: unit → int           bsp_g: unit → float           bsp_l: unit → float
mkpar: (int → α) → α par
apply: (α → β) par → α par → β par
type α option = None | Some of α
put: (int → α option) par → (int → α option) par       at: α par → int → α

```

Figure 2: The core BSMLlib library

where  $s_d^j$  is the number of super-steps of the BSP unit  $j$  during the  $d$ -step  $d$  and  $w_b^a$  the local processing time on processor  $b$  during each super-step  $a$  and  $h_b^{(a)} = \max\{h_{b+}^{(a)}, h_{b-}^{(a)}\}$  where  $h_{b+}^{(a)}$  (resp.  $h_{b-}^{(a)}$ ) is the number of words transmitted (resp. received) by processor  $b$  during each super-step  $a$ . The execution time for a program is thus bounded by:  $\Psi = \max\{\Phi_{R,j}/j \in \{0, 1, \dots, P-1\}\}$ , where  $R$  is the number of  $d$ -steps of the program. The DMM model takes into account that a BSP unit only synchronizes with its incoming partner and is therefore more accurate than the BSP<sup>2</sup> one: more algorithms for irregular problems could be analyzed efficiently.

### 3. The DMML Language

There is currently no implementation of a full DMML language but rather a partial implementation as a *library* for Objective Caml (OCaml) [11]. In this section we describe the BSMLlib core and the new primitives added to obtain DMML.

#### 3.1. The BSMLlib library

The core of so-called BSMLlib is based on the elements given in Figure 2. It gives access to the BSP parameters of the underlying architecture. In particular, **bsp\_p**() is  $p$ , the *static* number of processes. There is an abstract polymorphic type  $\alpha$  **par** which represents the type of  $p$ -wide *parallel vectors* of objects of type  $\alpha$  one per process. The nesting of **par** types is prohibited. Our *type system* enforces this restriction [9]. The BSML parallel constructs operate on parallel vectors. Those parallel vectors are created by **mkpar** so that **(mkpar f)** stores **(f i)** on process  $i$  for  $i$  between 0 and  $(p-1)$ . We usually write **f** as **(fun pid → e)** to show that the expression  $e$  may be different on each processor. This expression  $e$  is said to be *local*. The expression **(mkpar f)** is a parallel object and it is said to be *global*.

Asynchronous phases are programmed with **mkpar** and **apply**. The expression **(apply (mkpar f) (mkpar e))** stores **((f i)(e i))** on process  $i$ . The communication and synchronization phases are expressed by **put**. Consider the expression: **put(mkpar(fun i → fs<sub>i</sub>))**. To send a value  $v$  from process  $j$  to process  $i$ , the function  $fs_j$  at process  $j$  must be such that **(fs<sub>j</sub> i)** evaluates to **Some v**. To send no value from process  $j$  to process  $i$ , **(fs<sub>j</sub> i)** must evaluate to **None**. This expression evaluates to a parallel vector containing a function  $fd_i$  of delivered messages on every process. At process  $i$ , **(fd<sub>i</sub> j)** evaluates to **None** if process  $j$  sent no message to process  $i$  or evaluates to **Some v** if process  $j$  sent the value  $v$  to the process  $i$ .

**at** is the synchronous projection primitive where **(at vec n)** returns the  $n^{\text{th}}$  value of the parallel vector **vec**. **at** expresses communication and synchronization phases. Without it, the global control cannot take into account data computed locally. Global conditional is necessary of express algorithms like:

<b>dm_bsp_p</b>	: int → int	<b>dm_p</b>	: unit → int
<b>dm_bsp_s</b>	: int → float	<b>dm_g</b>	: unit → float
<b>dm_bsp_g</b>	: int → float	<b>dm_l</b>	: unit → float
<b>dm_bsp_l</b>	: int → float		
<b>mkdep</b>	: (int → $\alpha$ ) → $\alpha$ <b>dep</b>		
<b>applydep</b>	: ( $\alpha$ → $\beta$ ) <b>dep</b> → $\alpha$ <b>dep</b> → $\beta$ <b>dep</b>		
<b>get</b>	: (int → int → int option) <b>par dep</b> → (int → $\alpha$ option) <b>par dep</b> → (int → int → $\alpha$ option) <b>par dep</b>		
<b>atdep</b>	: $\alpha$ <b>par dep</b> → int → int → $\alpha$		

Figure 3: The added primitives

**Repeat** Parallel Iteration **Until** Max of local errors  $< \epsilon$

The projection should not be evaluated inside the scope of a **mkpar**. This is enforced by our type system [9]. The following program is a small example of a direct broadcast algorithm in BSML, where noSome is such as noSome (Some x)=x:

**exception** Bcast

```

let replicate x = mkpar (fun pid → x) and parfun f vv = apply (replicate f) vv
let bcast_direct rt vv = if (root<0 || root>=bsp_p()) then raise Bcast else
  let msg = mkpar (fun pid v dst → if pid=root then Some v else None)
  in parfun noSome (apply (put (apply msg vv)) (replicate root))

```

### 3.2. The DMML library

DMML extends BSMLlib by adding new primitives on a new level called departmental. The core of this library adds the primitives given in Figure 3. DMML offers functions to access to the parameters of the metacomputer, in particular, the function **dm\_p**:unit → int (resp. **dm\_g** and **dm\_l**) is such that the value of **dm\_p**() is  $P$ , the static number of BSP units (resp.  $G$  and  $L$ , bandwidth and latency of the departmental WAN). Parameters of the BSP units are available through the functions **dm\_bsp\_p**, **dm\_bsp\_s**, **dm\_bsp\_g** and **dm\_bsp\_l**. For example (**dm\_bsp\_p** a) gives the number of processors of the  $a^{\text{th}}$  parallel machine.

There is also a new polymorphic type  $\alpha$  **dep** which represents the type of  $P$ -wide departmental vectors of objects of type  $\alpha$  one per BSP unit. The nesting of **dep** into **dep** or into **par** types is prohibited. But the  $\alpha$  of a **dep** type could be either a usual OCaml value or a BSML value.

The DMML departmental constructs operate on departmental vectors. Those vectors are created by **mkdep** so that (**mkdep** f) stores (f a) on the BSP unit  $a$  for  $a$  between 0 and  $(P - 1)$ . The BSML parallel values should not be evaluated outside the scope of a **mkdep**. This could be enforced by a type system, but for the moment, the programmer is responsible for respecting this rule. The BSML parallel constructs operate on parallel vectors of size  $p_a$  for the BSP unit  $a$ . For example in the scope of a **mkdep**, (**mkpar** f) stores (f i) on process  $i$  for  $i$  between 0 and  $(p_a - 1)$  for the BSP unit  $a$ . The expression (**mkdep** f) is said to be *departmental*.

A DMM algorithm is expressed as a combination of asynchronous BSP computations (first phase of a  $d$ -step) and phases of communications (second phase of a  $d$ -step). Asynchronous phases are programmed with **mkdep** and with **applydep**. This primitives is such as on BSP unit  $a$ , (**applydep** (**mkdep** f) (**mkdep** e)) stores

((f a)(e a)). Communications are expressed by **get**. Consider the expression:

```
get (mkdep (fun a → mkpar (fun i → fa,i)))
    (mkdep (fun b → mkpar (fun j → vb,j)))
```

For a process  $i$  of the BSP unit  $a$ , to receive the  $n^{\text{th}}$  value from the process  $j$  of the BSP unit  $b$  (it is an incoming partners), the function  $f_{a,i}$  at process  $i$  of the BSP unit  $a$  must be such that  $(f_{a,i} b j)$  evaluates to **Some**  $n$ . To receive no value  $(f_{a,i} b j)$  must evaluate to **None**.

Our expression evaluates to a departmental vector containing parallel vectors of functions  $f'_{a,i}$  of delivered messages on every process of every BSP unit.

At process  $i$  of the BSP unit  $a$   $(f'_{a,i} b j)$  evaluates to **None** if process  $i$  of the BSP unit  $a$  receives no message from process  $j$  of the BSP unit  $b$  or if  $(v_{b,j} n)$  evaluates to **None** (the process  $j$  of the BSP unit  $b$  does not have a  $n^{\text{th}}$  value and sends the empty value). It also evaluates to **Some**  $v_{b,j}^n$  if it received a value from the process  $j$  of the BSP unit  $b$  and if  $(v_{b,j} n)$  evaluates to **(Some**  $v_{b,j}^n$ **)**.

There is also a projection primitive **atdep**. It is used in the same way as the **at** primitive but it takes as arguments a departmental vector of parallel vectors and two integers being the number of the cluster and the number of the process considered. This primitive should not be evaluated inside a **mkdep**. Using **atdep** the departmental behavior of the program could depend on a local value. The following program is a small example of a direct broadcast algorithm in DMML:

```
let replicate_all x = mkdep (fun a → replicate x)
let apply_all gf gv = applydep (applydep (mkdep (fun a f v → apply f v)) gf) gv
let parfun_all f x = apply_all (replicate_all f) x
```

```
let get_one_all datas srcs =
  let send=parfun_all (fun v n → Some v) datas
  and n_srcs=parfun_all (fun(a,i) → let ap = (natmod a (dm_p())) in
    (ap,(natmod i (dm_bsp_p ap)))) srcs in
  let ask = parfun_all (fun (a,i) cluster pid →
    if (cluster=a)&&(pid=i) then Some 0 else None) n_srcs in
  parfun2_all(fun f (a,i) → (noSome (f a i)))(get ask send) n_srcs
```

```
let bcast_direct_all rclus rpid vv =
  if rclus<0||rpid<0||rclus>=dm_p()||rpid>=(dm_bsp_p rclus)
  then raise Bcast
  else get_one_all vv (replicate_all(rclus,rpid))
```

Note that in the above example, the four last lines only are specific to the direct broadcast algorithm. The remaining functions are part of the DMML standard library and they could and are used in many other cases.

Thus this example gives a taste of the DMML programming style. The communications functions are purely functional – which is a high level feature, not so common and a requirement to use the extraction capability of a proof assistant – but the explicit handling of messages could make them difficult to use and could seem not so high level. Nevertheless it is very easy, and concise, to build more and more complex functions, in particular using the higher-order nature of functions. The user of the DMML library is offered a complete set of functions (maps,

broadcasts, folds, scans, *etc.*) implemented with the primitives only and can still implement its own functions using the primitives. DMML programs can use the DMM parameters to adapt themselves to the underlying architecture, which make them portable (in the same way the BSP algorithms are portable).

#### 4. Formal Semantics

Reasoning on the complete definition of a functional and parallel language such as DMML, would have been complex and tedious. In order to simplify the presentation and to ease the formal reasoning, this section introduces a *core* language.

##### 4.1. Syntax of the DMML core-language

*Expressions*, written  $e$  and variants, have the following abstract syntax:

$$\begin{aligned}
 e ::= & x \mid c \mid op \mid \mathbf{fun} \ x \rightarrow e \mid (e \ e) \mid \mathbf{if} \ e \ \mathbf{then} \ e \ \mathbf{else} \ e \mid \mathbf{let} \ x = e \ \mathbf{in} \ e \\
 & \mid \mathbf{mkpar} \ e \mid \mathbf{apply} \ e \ e \mid \mathbf{put} \ e \mid \mathbf{at} \ e \ e \\
 & \mid \mathbf{mkdep} \ e \mid \mathbf{applydep} \ e \ e \mid \mathbf{get} \ e \ e \mid \mathbf{atdep} \ e \ e \ e
 \end{aligned}$$

In this grammar,  $x$  ranges over a countable set of identifiers. Constants  $c$  are  $()$  (the only value of type `unit` in OCaml), the integers, the booleans and the value `nc` (which stands for no communication) which plays the role of the `None` constructor in OCaml for the `put` communication primitive. The set of predefined operations  $op$  contains arithmetic and boolean operations, the test function `isnc` of the `nc` constant, the `fix` operator and the function of access to the parameters of the metacomputer.

This syntax is the programmer's syntax. Reduction can produce additional expressions, enumerated parallel vectors and enumerated departmental vectors:

$$e ::= \dots \mid \langle e, \dots, e, \dots, e \rangle \mid \llbracket e, \dots, e, \dots, e \rrbracket$$

There is one semantics per size of the DMM machine. We will note  $a, b$  and variants, the BSP units, and  $i, j$  and variants the processors. By size we mean the number of BSP units, and  $p_a$  the number of processors of the BSP unit  $a$ , for each unit of the metacomputer. The size of departmental vectors is  $P$  and the sizes of parallel vectors can be any  $p_a$ .

Before presenting the dynamic semantics of our core-language, i.e, how the expressions are computed to values, we present the values themselves:

$$v ::= \mathbf{fun} \ x \rightarrow e \mid c \mid op \mid \langle v, \dots, v, \dots, v \rangle \mid \llbracket v, \dots, v, \dots, v \rrbracket$$

##### 4.2. Reduction rules

To express the DMML semantics, we use a *small-step* semantics. It consists of a predicate between an expression and another expression. The small-step semantics describes all the steps from an expression to a value. The small-steps semantics has the following form:  $e \rightarrow e'$ . We note  $\rightarrow^*$ , for the reflexive transitive closure of  $\rightarrow$ , e.g., we note  $e_0 \rightarrow^* v$  for  $e_0 \rightarrow e_1 \rightarrow e_2 \rightarrow \dots \rightarrow v$ .

To define the relation  $\rightarrow$ , we begin to define three relations, one for each kind of expression: local (usual OCaml expressions), parallel (BSML expressions) and departmental (DMML expressions).

$$(\mathbf{fun} \ x \rightarrow e) \ v \xrightarrow{\varepsilon} e[x \leftarrow v] \qquad (\mathbf{let} \ x = v \ \mathbf{in} \ e) \xrightarrow{\varepsilon} e[x \leftarrow v]$$

Figure 4: Head reductions

$$\begin{array}{ll}
 \mathbf{fix}(\mathbf{fun} \ x \rightarrow e) \xrightarrow{\varepsilon} e[x \leftarrow \mathbf{fix}(\mathbf{fun} \ x \rightarrow e)] & \mathbf{fix}(\mathbf{op}) \xrightarrow{\varepsilon} \mathbf{op} \\
 (\mathbf{isnc} \ v) \xrightarrow{\varepsilon} \mathbf{false} \ \text{if } v \neq \mathbf{nc} & (\mathbf{isnc} \ \mathbf{nc}) \xrightarrow{\varepsilon} \mathbf{true} \\
 \mathbf{if} \ \mathbf{false} \ \mathbf{then} \ e_1 \ \mathbf{else} \ e_2 \xrightarrow{\varepsilon} e_2 & (\mathbf{dm\_bsp\_p} \ a) \xrightarrow{\varepsilon} \mathbf{P}_a \\
 \mathbf{if} \ \mathbf{true} \ \mathbf{then} \ e_1 \ \mathbf{else} \ e_2 \xrightarrow{\varepsilon} e_1 & (\mathbf{dm\_p} \ ()) \xrightarrow{\varepsilon} \mathbf{P} \\
 \mathbf{access} \ [v_0, \dots, v_n, \dots, v_i] \ n \xrightarrow{\varepsilon} v_n & \mathbf{init} \ l \ f \xrightarrow{\varepsilon} [(f \ 0), \dots, (f \ (l-1))]
 \end{array}$$

Figure 5: Predefined operations

These three relations are :

- $e \xrightarrow{i,a} e'$ : at process  $i$  of the BSP unit  $a$  the expression  $e$  is reduced to  $e'$ ;
- $e \xrightarrow{\mathbb{M}_a} e'$ : at parallel machine  $a$ , the expression  $e$  is reduced to  $e'$ ;
- $e \xrightarrow{\approx} e'$ : the expression  $e$  is reduced to  $e'$  by the whole metacomputer.

Each kind of expression, local, parallel or departmental contains usual function abstractions and applications. Thus all these relations contain the following relation  $\xrightarrow{\varepsilon}$ , called the relation of *head reduction* given in figure 4

We write  $e_1[x \leftarrow e_2]$  the expression obtained by substituting all the free occurrences of  $x$  in  $e_1$  by  $e_2$ . Free occurrences of a variable is defined as a classical and trivial inductive function on our expressions.

Some rules, the  $\delta$ -rules, for predefined operations are given in figure 5. Rules for parallel and departmental primitives and given respectively in figures 6 and 7.

We give here the semantics of the BSML primitive **put** in two steps, corresponding accurately to the current implementation of BSML. Firstly each processor creates a purely functional array of values by applying the function which it holds to all the possible numbers of processor in the unit where the reduction takes place. Then a lower level primitive **send** makes the exchanges and returns a parallel vector of arrays. The value at index  $j$  of the array is sent to the processor  $j$  if it is not **nc**. Processor  $j$  receives it and stores it at index  $i$  of result array. The function **mkf** builds the parallel vector of result functions from the parallel vector of arrays.

To do that we introduce new expressions : arrays, written  $[e, \dots, e, \dots, e]$ . The two predefined operations **init** and **access** operates on arrays (figure 5).

Rules of primitives **mkdep**, **applydep** and **atdep** are similar to the rules for BSML primitives but differ by the number of arguments and by kinds of vectors.

**get** needs, like **put**, two lower level primitives: **senddep** and **mkanswer**.

The first primitive (rule 10) adapts at the departmental level the **send** of the global level. The argument is a departmental vector of parallel vectors of arrays of arrays of values. The results is a value which has the "same type". The argument of this operation is such as at a given processor  $i$  of a unit  $a$  the array of arrays indicates for each pair  $(b, j)$  of unit and processor, the value to send to processor  $j$  of unit  $b$ . The result is such as each processor  $i$  of a unit  $a$  the array of arrays

$$\mathbf{mkpar} \ v \xrightarrow[\mathfrak{K}_a]{\varepsilon} \langle (v \ 0), \dots, (v \ (p_a - 1)) \rangle \quad (1)$$

$$\mathbf{apply} \ \langle v_0, \dots, v_{p_a-1} \rangle \ \langle v'_0, \dots, v'_{p_a-1} \rangle \xrightarrow[\mathfrak{K}_a]{\varepsilon} \langle (v_0 \ v'_0), \dots, (v_{p_a-1} \ v'_{p_a-1}) \rangle \quad (2)$$

$$\mathbf{at} \ \langle \dots, v_n, \dots \rangle \ n \xrightarrow[\mathfrak{K}_a]{\varepsilon} \ v_n \quad (3)$$

$$\mathbf{put} \ \langle v_0, \dots, v_i, \dots, v_{p_a-1} \rangle \xrightarrow[\mathfrak{K}_a]{\varepsilon} \ (\mathbf{mkf} \ (\mathbf{send} \ \langle \dots, (\mathbf{init} \ p_a \ v_i), \dots \rangle)) \quad (4)$$

$$\mathbf{send} \ \langle [v_0^0, \dots, v_0^{p_a-1}], \dots, [v_{p_a-1}^0, \dots, v_{p_a-1}^{p_a-1}] \rangle \xrightarrow[\mathfrak{K}_a]{\varepsilon} \ \langle [v_0^0, \dots, v_{p_a-1}^0], \dots, [v_0^{p_a-1}, \dots, v_{p_a-1}^{p_a-1}] \rangle \quad (5)$$

where  $\mathbf{mkf} = \mathbf{apply} \ (\mathbf{mkpar} \ (\mathbf{fun} \ a \ t \ i \ \rightarrow$   
 $\mathbf{if} \ (0 \leq i) \ \& \ (i < (\mathbf{dm\_bsp\_p} \ a)) \ \mathbf{then} \ (\mathbf{access} \ t \ i) \ \mathbf{else} \ \mathbf{nc}))$

Figure 6: Reduction  $\xrightarrow[\mathfrak{K}_a]{\varepsilon}$

indicates for each pair  $(b, j)$  of unit and processor, the value the processor  $j$  of unit  $b$  sent to processor  $i$  of unit  $a$ .

The second primitive (rule 11) takes as argument a departmental vector of parallel vectors of arrays of arrays of integers and value  $\mathbf{nc}$  which indicates to the processor  $i$  of unit  $a$ , at the index  $j$  in array with the index  $b$  in the arrays of arrays, the number of the value requested by the processor  $j$  of the unit  $b$  to the processor  $i$  of the unit  $a$ . The values are given by the departmental vector of parallel vectors of functions, second argument of this primitive. The result is a departmental vector of parallel vectors of arrays of arrays of the values to be transmitted.

The rule 9 allows the following steps: 1) A departmental vector of parallel vectors of arrays of arrays indicating which numbers of values must be requested is created from the first argument of the  $\mathbf{get}$ . It is the reduction of the expressions  $t_i^a$  which will create these arrays of arrays. 2) These numbers of values are transmitted to the concerned processors: it is the primitive  $\mathbf{senddep}$  applied to the departmental vector describes at the preceding step which carries out this. 3) The values to be sent are prepared using the  $\mathbf{mkanswer}$  primitive. 4) The values are sent with the second call to  $\mathbf{senddep}$ . 5) Then the function  $\mathbf{mkf2}$  transforms the departmental vector of parallel vectors of arrays of arrays into a departmental vector of parallel vectors of functions, result of the  $\mathbf{get}$ .

The three relations are then obtained as follows:

$$\xrightarrow[\mathfrak{K}_a]{\varepsilon} = \xrightarrow[\delta]{\varepsilon} \cup \xrightarrow[\delta]{\varepsilon} \quad \text{and} \quad \xrightarrow[\mathfrak{K}_a]{\varepsilon} = \xrightarrow[\mathfrak{K}_a]{\varepsilon} \cup \xrightarrow[\delta]{\varepsilon} \cup \xrightarrow[\delta]{\varepsilon} \quad \text{and} \quad \xrightarrow[\diamond]{\varepsilon} = \xrightarrow[\diamond]{\varepsilon} \cup \xrightarrow[\delta]{\varepsilon} \cup \xrightarrow[\delta]{\varepsilon}$$

### 4.3. Context rules

It is easy to see that we cannot always make a head reduction: we have to reduct within an expression. To define this *deep reduction*, we define some kind of *contexts*, i.e., an expression with a *hole* noted  $\square$  that have the abstract syntax given in Figure 8. The hole gives where expressions could be reduced. In this way, the contexts give the order of evaluation of the arguments of the constructions of the language, i.e., the *strategie*. We note  $\mathbf{op}$  is a parallel or departmental primitive. The  $\Gamma$  context is used to define a *departmental reduction* of the metacomputer, i.e., a reduction

$$\mathbf{mkdep} \ v \xrightarrow[\diamond]{\varepsilon} \langle\langle (v \ 0), \dots, (v \ (P-1)) \rangle\rangle \quad (6)$$

$$\mathbf{applydep} \ \langle\langle v_0, \dots, v_{P-1} \rangle\rangle \xrightarrow[\diamond]{\varepsilon} \langle\langle (v_0 \ v'_0), \dots, (v_{P-1} \ v'_{P-1}) \rangle\rangle \quad (7)$$

$$\mathbf{atdep} \ \langle\langle \dots, \langle \dots, v_i^a, \dots \rangle, \dots \rangle\rangle \ a \ i \xrightarrow[\diamond]{\varepsilon} v_i^a \quad (8)$$

$$\mathbf{get} \ \langle\langle \dots, \langle \dots, f_i^a \dots \rangle, \dots \rangle\rangle \xrightarrow[\diamond]{\varepsilon} \langle\langle \dots, \langle \dots, g_i^a \dots \rangle, \dots \rangle\rangle \quad (9)$$

$$\begin{aligned} & (\mathbf{mkf2} \ (\mathbf{senddep} \ (\mathbf{mkanswer} \\ & \quad (\mathbf{senddep} \ \langle\langle \dots, \langle \dots, t_i^a \dots \rangle, \dots \rangle\rangle) \\ & \quad \langle\langle \dots, \langle \dots, g_i^a \dots \rangle, \dots \rangle\rangle)) \ \text{where} \\ & \quad t_i^a = (\mathbf{init} \ P \ (\mathbf{fun} \ b \rightarrow \\ & \quad \quad (\mathbf{init} \ (\mathbf{dm\_bsp\_p} \ b) \ (f_i^a \ b)))) \end{aligned}$$

$$\mathbf{senddep} \ \langle\langle \dots, \langle \dots, t_i^a, \dots \rangle, \dots \rangle\rangle \ \text{where} \quad \langle\langle \dots, \langle \dots, t_i^a, \dots \rangle, \dots \rangle\rangle \ \text{where} \quad (10)$$

$$t_i^a = \begin{bmatrix} [n_{(i,0)}^{(a,0)}, \dots, n_{(i,p_0-1)}^{(a,0)}], \\ \dots, \\ [n_{(i,0)}^{(a,P-1)}, \dots, n_{(i,p_{P-1}-1)}^{(a,P-1)}] \end{bmatrix} \quad t_i^a = \begin{bmatrix} [n_{(0,i)}^{(0,a)}, \dots, n_{(p_0-1,i)}^{(0,a)}], \\ \dots, \\ [n_{(0,i)}^{(P-1,a)}, \dots, n_{(p_{P-1}-1,i)}^{(P-1,a)}] \end{bmatrix}$$

$$\mathbf{mkanswer} \ \langle\langle \dots, \langle \dots, t_i^a \dots \rangle, \dots \rangle\rangle \xrightarrow[\diamond]{\varepsilon} \langle\langle \dots, \langle \dots, g_i^a \dots \rangle, \dots \rangle\rangle \quad (11)$$

$$\begin{aligned} & \langle\langle \dots, \langle \dots, t_i^a \dots \rangle, \dots \rangle\rangle \ \text{where} \\ & \quad (\mathbf{init} \ P \ (\mathbf{fun} \ b \rightarrow (\mathbf{init} \ (\mathbf{dm\_bsp\_p} \ b) \\ & \quad \quad (\mathbf{fun} \ j \rightarrow \mathbf{let} \ v = \\ & \quad \quad \quad (\mathbf{access} \ (\mathbf{access} \ t_i^a \ b) \ j) \ \mathbf{in} \\ & \quad \quad \quad \mathbf{if} \ (\mathbf{isnc} \ v) \ \mathbf{then} \ \mathbf{nc} \ \mathbf{else} \ (g_i^a \ v)))))) \end{aligned}$$

$$\text{where } \mathbf{mkf2} = \begin{cases} \mathbf{applydep} \ (\mathbf{mkpdep} \ (\mathbf{fun} \ a \rightarrow \\ \quad \mathbf{apply} \ (\mathbf{mkpar} \ (\mathbf{fun} \ i \ t \ b \ j \rightarrow \\ \quad \quad \mathbf{if} \ ((0 \leq b) \ \& \ (b < P) \ \& \ (0 \leq j) \ \& \ (j < (\mathbf{dm\_bsp\_p} \ b))) \\ \quad \quad \mathbf{then} \ (\mathbf{access} \ (\mathbf{access} \ t \ b) \ j) \ \mathbf{else} \ \mathbf{nc}))) \end{cases}$$

 Figure 7: Reduction  $\xrightarrow[\diamond]{\varepsilon}$ 

outside a departmental vector. For example:

$$\Gamma = \mathbf{let} \ x = \square \ \mathbf{in} \ \mathbf{mkdep} \ (\mathbf{fun} \ clus \rightarrow \dots)$$

The reduction will occur at the hole to first compute the value of  $x$ . The  $\Gamma^a$  context is used to define in which component of a departmental vector the reduction is done, i.e., which BSP unit  $a$  reduces its global expression. This context uses the  $\Gamma_g$  context which defines a *global reduction* on a BSP unit. Note that, in this way, the hole is inside a departmental vector. For example, the following context:  $\Gamma^a = \mathbf{applydep} \ v \ \langle\langle v_0, e_1, \dots, \Gamma_g \rangle\rangle$  and  $\Gamma_g = \mathbf{mkpar} \ \square$  is used to define that the last BSP unit first computes the argument of the  $\mathbf{mkpar}$  primitive.

The  $\Gamma_i^a$  context is used to define a *local reduction* at processor  $i$  on a parallel machine  $a$ : first the context finds a hole in a departmental vector, next, a  $\Gamma_i$  context finds a hole in a parallel vector, i.e., which processor makes the reduction, and to end a  $\Gamma_l$  context finds the hole in a local expression, i.e., standard OCaml expression. Note that this hole is inside a parallel vector which is inside a departmental vector.

Now we can reduce in-depth in the sub-expressions. To define this deep reduction, we use the inference rules of all the different kinds of context rules:

$$\begin{array}{l}
 \Gamma_i^a ::= \Gamma_i^a e \\
 \quad | v \Gamma_i^a \\
 \quad | \text{let } x = \Gamma_i^a \text{ in } e \\
 \quad | \text{if } \Gamma_i^a \text{ then } e \text{ else } e \\
 \quad | \text{op } \Gamma_i^a e_2 \dots e_n \\
 \quad | \text{op } v_1 \Gamma_i^a \dots e_n \\
 \quad | \dots \\
 \quad | \text{op } v_1 v_2 \dots \Gamma_i^a \\
 \quad | \langle \langle e, \dots, \overbrace{\Gamma_i^a}^a, e, \dots, e \rangle \rangle
 \end{array}
 \quad
 \begin{array}{l}
 \Gamma^a ::= \Gamma^a e \\
 \quad | v \Gamma^a \\
 \quad | \text{let } x = \Gamma^a \text{ in } e \\
 \quad | \text{if } \Gamma^a \text{ then } e \text{ else } e \\
 \quad | \text{op } \Gamma^a e_2 \dots e_n \\
 \quad | \text{op } v_1 \Gamma^a \dots e_n \\
 \quad | \dots \\
 \quad | \text{op } v_1 v_2 \dots \Gamma^a \\
 \quad | \langle \langle e, \dots, \overbrace{\Gamma^a}^a, e, \dots, e \rangle \rangle
 \end{array}
 \quad
 \begin{array}{l}
 \Gamma ::= [] \\
 \quad | \Gamma e \\
 \quad | v \Gamma \\
 \quad | \text{let } x = \Gamma \text{ in } e \\
 \quad | \text{if } \Gamma \text{ then } e \text{ else } e \\
 \quad | \text{op } \Gamma e_2 \dots e_n \\
 \quad | \text{op } v_1 \Gamma \dots e_n \\
 \quad | \dots \\
 \quad | \text{op } v_1 v_2 \dots \Gamma
 \end{array}$$
  

$$\begin{array}{l}
 \Gamma_i ::= \Gamma_i e \\
 \quad | v \Gamma_i \\
 \quad | \text{let } x = \Gamma_i \text{ in } e \\
 \quad | \text{if } \Gamma_i \text{ then } e \text{ else } e \\
 \quad | \text{op } \Gamma_i e_2 \dots e_n \\
 \quad | \text{op } v_1 \Gamma_i \dots e_n \\
 \quad | \dots \\
 \quad | \text{op } v_1 v_2 \dots \Gamma_i \\
 \quad | \langle e, \dots, \overbrace{\Gamma_i}^i, e, \dots, e \rangle
 \end{array}
 \quad
 \begin{array}{l}
 \Gamma_g ::= [] \\
 \quad | \Gamma_g e \\
 \quad | v \Gamma_g \\
 \quad | \text{let } x = \Gamma_g \text{ in } e \\
 \quad | \text{if } \Gamma_g \text{ then } e \text{ else } e \\
 \quad | \text{op } \Gamma_g e_2 \dots e_n \\
 \quad | \text{op } v_1 \Gamma_g \dots e_n \\
 \quad | \dots \\
 \quad | \text{op } v_1 v_2 \dots \Gamma_g
 \end{array}
 \quad
 \begin{array}{l}
 \Gamma_l ::= [] \\
 \quad | \Gamma_l e \\
 \quad | v \Gamma_l \\
 \quad | \text{let } x = \Gamma_l \text{ in } e \\
 \quad | \text{if } \Gamma_l \text{ then } e \text{ else } e \\
 \quad | [\Gamma_l, e_2, \dots, e_n] \\
 \quad | [v_0, \Gamma_l, \dots, e_n] \\
 \quad | \dots \\
 \quad | [v_0, v_1, \dots, \Gamma_l]
 \end{array}$$

Figure 8: Contexts

$$\frac{e \xrightarrow{i,a} e'}{\Gamma_i^a[e] \xrightarrow{i,a} \Gamma_i^a[e']}
 \quad
 \frac{e \xrightarrow{\kappa_a} e'}{\Gamma^a[e] \xrightarrow{\kappa_a} \Gamma^a[e']}
 \quad
 \frac{e \xrightarrow{\cong} e'}{\Gamma[e] \xrightarrow{\cong} \Gamma[e']}$$

**Proposition 1** *Let  $e$  be an expression. If  $e \xrightarrow{*} v_1$  and  $e \xrightarrow{*} v_2$  then  $v_1 = v_2$ .*

**Sketch of the Proof.** All the  $\delta$ -rules and head reductions, i.e., the axioms, are deterministic (local, global and departmental ones). The rules are not always deterministic, i.e., several axioms can be applied at the same time, parallelism comes from the context rules. But if a context gives two possible reductions in a parallel or departmental vector, it is easy to see that these two reductions could be done in any order and give the same result because a reduction does not affect the result of the other one. In this way the DMML language is confluent.  $\square$

## 5. Examples and Implementation

### 5.1. Broadcast

In the broadcast program, a single process of a BSP unit, called the root  $r$ , sends a message to all other processes. It could be done in a direct way: each process of each BSP unit asks the value of the root (see example of section 3) and the cost is:

$$\max \left\{ \begin{array}{l} a \in \{0 \dots P-1\} \setminus \{r\} \quad (\mathcal{S}(v) \times p_a \times (g_a + G + g_r) + l_r + l_a + L) \\ \text{and } \sum_{i \in \{0 \dots P-1\}} (g_r \times p_i \times \mathcal{S}(v)) + l_r \end{array} \right.$$

where  $v$  is the sent value and  $\mathcal{S}$  the size in bytes of the value. The cost of the program is the maximum time for a BSP unit to receive the value and for the root to send it to all the processes of all the BSP units.

Another way is that each process of each BSP unit receives from the root only a subpart of the message. Each BSP unit contains all the parts needed to rebuild the initial value. Then on each BSP unit there is a total exchange of these parts to

obtain the whole message. Thus we have the following cost:

$$\max \left\{ \begin{array}{l} a \in \{0 \dots P-1\} \setminus \{r\} \quad (\mathcal{S}(v) \times (g_a + G + g_r) + l_r + l_a + L + (p_a \times g_a \times \lceil \frac{\mathcal{S}(v)}{p_a} \rceil + l_a)) \\ \text{and } g_r \times P \times \mathcal{S}(v) + l_r + (p_a \times g_r \times \lceil \frac{\mathcal{S}(v)}{p_r} \rceil + l_r) \end{array} \right.$$

The cost of the program is the maximum time for a BSP unit to receive the parts of the value and to totally exchange them and for the root to send the parts and to totally exchange them. Note that this program uses a program to scatter the message from the root.

### 5.2. Departmental Reduction

Our second example is the classical parallel reduction: each process of each BSP unit contains a value and we want to obtain the sum of these values. For this, a naive algorithm could exchange all the needed values and then each process performs a local reduction. In this way, the cost of this program is close to the BSP cost of the direct algorithm.

As an example of a formal cost analysis (we refer to a report available at [dmmlib.free.fr](http://dmmlib.free.fr) for more details about how to formally give costs to DMML programs) of a less naive algorithm, we choose the multiplication-reduction of polynomials: each process contains a polynomial and we want to compute their global multiplication. We make the following hypotheses: 1) The clusters are sorted by their efficiency to perform a BSP reduction 2) the coefficients of the polynomials ( $\sum_{i=0}^m c_i X^i$ ) are stored in an array of floats such that  $c_i$  is located at position  $i$ . We write  $\mathcal{S}(n)$  for the size of a polynomial of degree  $n$ . In this way, we have the following property:  $\mathcal{S}(poly1 \times poly2) = \mathcal{S}(poly1) + \mathcal{S}(poly2)$  if we make the hypothesis that the size of a float does not depend of its value.

The algorithm runs as follow. First each BSP unit  $a$  performs a direct BSP reduction. The cost for each of them is thus:

$$A_a = (n \times p_a)^2 \times r_a + (p_a - 1) \times \mathcal{S}(n) \times g + l_a$$

where  $n$  is the maximal degree of the polynomials and  $r_a$  the time to perform a float multiplication. Second the root process of each BSP units receives the polynomials of the previous BSP unit. In this way, the cost to receive the polynomials is:

$$B_a = \sum_{\forall b < a} (l_b + (g_b + G + g_a) \times \mathcal{S}(p_b \times n) + l_a + L)$$

and to send them is:  $C_a = l_a + \mathcal{S}(p_a \times n) \times g_a \times (P - a)$ .

With these received polynomials, the BSP unit is able to finish its reduction and the cost is  $D_a = (p_a - 1) \times (\sum_{\forall b < a} \mathcal{S}(p_b \times n) \times g_a + l_b + (\sum_{\forall b < a} \mathcal{S}(p_b \times n) \times r_a)$ . The execution time for the program is thus:  $\max_{a \in \{0 \dots P-1\}} (A_a + B_a + C_a + D_a)$ .

### 5.3. Implementation of the DMML library

There are two main versions of the DMML library: a sequential version and a parallel one. Based on a confluent semantics, the evaluation of a pure functional parallel program will lead to the same value with both versions. In the sequential version, parallel and departmental vectors are implemented with OCaml arrays.

In the parallel version, our primitives are implemented as SPMD programs. A parallel (resp. departmental) vector is supposed to contain one value per process

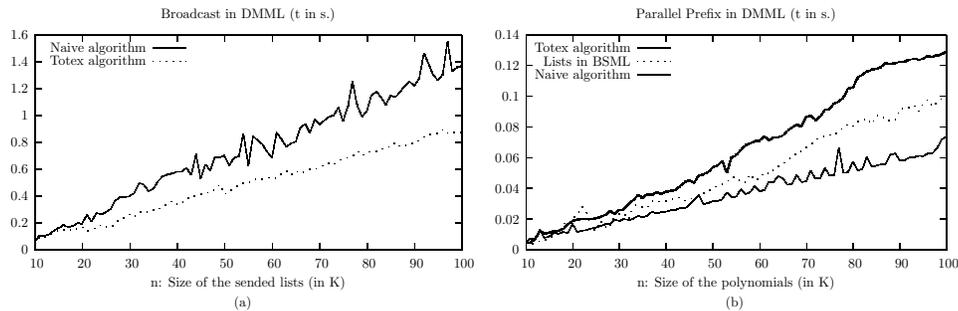


Figure 9: Benchmarks of collective operations

(resp. per BSP unit). The non-communicating primitives are thus very simple to implement using the “pid” of each BSP unit and each process.

Currently for the BSP part of our language, any MPI library can be used. In fact we only use a very small subpart of MPI: functions given the process identifier and the number of processes of the BSP unit and the all-to-all collective functions.

For the departmental part of our language, we use the thread facilities of the OCaml language: communication environments [12] of each process are needed to save the functional value of its  $d$ -steps (each process as a variable which count the  $d$ -step) and is thus implemented as a thread. The asynchronous **request** and **return** of the **get** primitive are also implemented as threads and use the TCP/IP facilities of the OCaml language to communicate the values. The examples and a first implementation of DMML are available at [dmmlib.free.fr](http://dmmlib.free.fr).

#### 5.4. Benchmarks

Preliminary experiments have been done on a metacomputer with 6 Pentium IV nodes cluster interconnected with a Gigabit Ethernet network and with 3 Celeron III nodes cluster interconnected with a Fast Ethernet network. The two clusters are interconnected with a slow Ethernet network. Figure 8 summarizes the timings. These programs were run 10 times and the average was taken. The naive broadcast algorithm is clearly slower than the second algorithm. Preliminary experiments of parallel reduction of multiplication of polynomials have been done to show a performance comparison between a BSP algorithm on the metacomputer and DMM algorithms. The BSML program has been only run on the first cluster and contain the same number of polynomials: 3 processes contains 2 polynomials. Using a second cluster and a less naive algorithm achieved a scalability improvement.

## 6. Related Work

The asynchronous nature of some parallel patterns like farms and pipelines or divide-and-conquer parallel algorithms hampers their efficient implementation with flat data parallel languages with global barriers. To overcome these limitations, the BSP PUB library [4] provides the capacity to partition the current BSP machine into several subsets. Several models [6] allowing *subset synchronization* have been proposed. The authors of the BSP Worldwide Standard Library report claims that an unwanted consequence of group partitioning is a loss of accuracy of the associated

performance model.

Another feature of PUB is the oblivious synchronization. It implies that different processors can be in different super-steps at the same time and thus the MPM model of [3] seems to be more adequate for a cost analysis of this kinds of programs. In our computation model, the two-tiered parallel levels are not based on subset synchronizations: it is a flat two-tiered level model and our cost model benefits of the advantages of BSP and MPM ones.

[20] presents an hierarchical extension of the BSP model with heterogeneous processors. But in this model, the execution of a program also proceeds in hyper-steps. Furthermore, the gateway is used for computation. The authors only analyze two-tiered level programs and have the same problems as in the BSP<sup>2</sup> model: the time of the global barrier of synchronization of a hyper-step. [5] is another hierarchical model with heterogeneous processors and asynchronous steps. The large number of parameters in this model introduce a hardly tractable complexity. The same problems occur in the model of [16]. Interesting work is the model of [2] for hierarchical and heterogeneous computers. But the main problem in this model is that the programs are difficult to analyze because the end-to-end bandwidth is combined with the latency. Moreover, in all those frameworks an execution model as in the BSP one lacks and deadlocks are possible.

Similar work to ours was conducted by [15], who performed an empirical study of the benefits of using a two-tiered parallel programming model. Their approach is based on data-duplication, all-to-all broadcasting and multicast message passing whereby these data would only be sent once between BSP units and then copied to all the BSP processors within the destination unit. In this way, they achieved a considerable scalability improvement. Another similar model is the pLogP model, parametrized Log-P of [10]. The authors introduce a two-tiered extension of the Log-GP model to optimize with the help of a cost analysis the collective operations of their own MPI library. But the authors do not present any formal semantics nor formal cost model and they use a low level language. To our knowledge, the DMML language is the first functional language for metacomputing with a formal semantics and a cost model.

## **7. Conclusions and Future Work**

Earlier research has shown that many parallel applications can be optimized to run efficiently on hierarchical wide-area systems. The BSP model has proved to be a trusty and worthy tool in the discipline of parallel programming for producing reliable and portable codes with predictable efficiency. However, additional complexity introduced by metacomputing forces a review of the model. We have considered a hierarchical extension of the BSP model, called the DMM model and we have also described a new functional parallel language for this new model. This language is based on a formal confluent semantics and allows programs cost analysis.

The first direction for future work is the design of algorithms for the DMM model and their implementations using the DMML library. To validate the cost model, we need a benchmark suite to determine the parameters of the metacomputer: this is an ongoing work. A complementary direction is to implement versions of DMML with adequate low level libraries for metacomputing [10,1]. Another direction of research

is the design of a distributed semantics of DMML, closer to the implementation, and the prove of its correctness with respect to the semantics presented here.

**Acknowledgements** This work is supported by a grant from the French Ministry of Research and the ACI Grid program (project CARAML). The authors wish to thank the anonymous referees for their comments. This paper is a revised version of the paper presented at CMPP'04.

## References

- [1] O. Aumage and al. Madeleine II: A Portable and Efficient Communication Library for High-Performance Cluster Computing. *Parallel Computing*, 28(4):607–626, 2002.
- [2] P.B. Bhat and al. Adaptative communication algorithms for distributed heterogeneous systems. *Parallel and Distributed Computation*, 59:252–279, 1999.
- [3] V. Blanco, J. A. González, C. León, C. Rodríguez, and G. Rodríguez. Predicting the performance of parallel programs. *Parallel Computing*, 30:337–356, 2004.
- [4] O. Bonorder, B. Juurlink, I. von Otte, and I. Rieping. The Paderborn University BSP (PUB) Library. Web pages at <http://www.uni-paderborn.de/~bsp/>, 2004.
- [5] F. Cappello, A.L. Rosenberg, and al. HiHCoHP toward a realistic communication model for hierarchical hyperclusters of heterogeneous processors. In *IPDPS'2001*. IEEE Press, 2001.
- [6] H. Cha and D. Lee. H-BSP: a Hierarchical BSP Computation Model. *Supercomputing*, 18(1):179–200, 2001.
- [7] F. Gava. Formal Proofs of Functional BSP Programs. *PPL*, 13(3):365–376, 2003.
- [8] F. Gava. Design of Departmental Metacomputing ML. In M. Bubak and al., editors, *ICCS 2004*, number 3038 in LNCS, pages 50–53. Springer Verlag, 2004.
- [9] F. Gava and F. Loulergue. A Static Analysis for Bulk Synchronous Parallel ML to Avoid Parallel Nesting. *FGCS*, 2004.
- [10] T. Kielmann, H. E. Bal, S. Gorchatch, K. Verstoep, and R. F. H. Hofman. Network performance-aware collective communication for clustered wide-area systems. *Parallel Computing*, 27:1431–1456, 2001.
- [11] X. Leroy, D. Doligez, J. Garrigue, D. Rémy, and J. Vouillon. The Objective Caml System release 3.08.1. Web pages at [www.ocaml.org](http://www.ocaml.org), 2004.
- [12] F. Loulergue, F. Gava, M. Arapinis, and F. Dabrowski. Semantics and Implementation of MSPML. *Int. Journal of Computer and Information Science*, 5(3), 2004.
- [13] F. Loulergue, G. Hains, and C. Foisy. A Calculus of Functional BSP Programs. *Science of Computer Programming*, 37(1-3):253–277, 2000.
- [14] J. M. R. Martin and A. Tiskin. BSP modelling a two-tiered parallel architectures. In B. M. Cook, editor, *WoTUG'99*, pages 47–55, 1999.
- [15] A. Plaatz, H. E. Bal, and al. Sensitivity of Parallel Applications to large differences in bandwidth and latency in two-layer interconnects. *FGCS*, 2004.
- [16] A.L. Rosenberg. Optimal sharing of partitionable workloads in heterogeneous networks of workstations. In *PDPTA '2000*, pages 413–419. CSREA Press, 2000.
- [17] D. B. Skillicorn, J. M. D. Hill, and W. F. McColl. Questions and Answers about BSP. *Scientific Programming*, 6(3):249–274, 1997.
- [18] L. Smarr and C. E. Catlett. Metacomputing. *CACM*, 35(6):44–52, 1992.
- [19] L.G. Valiant. A bridging model for parallel computation. *CACM*, 33(8):103, 1990.
- [20] T. L. Williams and R. J. Parsons. Exploiting hierarchy in heterogeneous environments. In *IEEE/ACM IPDPS'2001*, pages 140–147. IEEE Press, 2001.