

UNIVERSITÀ DI PISA
DIPARTIMENTO DI INFORMATICA

TECHNICAL REPORT: TR-07-05

Management in distributed systems: a semi-formal approach

Marco Aldinucci[◊] Marco Danelutto[◊] Peter Kilpatrick[•]

[◊]Dept. of Computer Science, University of Pisa

[•]Dept. of Computer Science, Queen's University Belfast

ADDRESS: via F. Buonarroti 2, 56127 Pisa, Italy. TEL: +39 050 2212700 FAX: +39 050 2212726

Management in distributed systems: a semi-formal approach ^{*}

Marco Aldinucci[◊] Marco Danelutto[◊] Peter Kilpatrick[•]

[◊]Dept. of Computer Science, University of Pisa

[•]Dept. of Computer Science, Queen's University Belfast

Abstract

Formal tools can be used in a “semi-formal” way to support distributed program analysis and tuning. We show how ORC has been used to reverse engineer a skeleton based programming environment and to remove one of the skeleton system’s recognized weak points. The semi-formal approach adopted allowed these steps to be performed in a programmer-friendly way.

Keywords. Distributed systems, orchestration, algorithmic skeletons.

1 Introduction

The *muskel* system, introduced by Danelutto in [1] and further elaborated in [2], reflects two modern trends in distributed system programming: the use of program skeletons and the provision of means for marshalling resources in the presence of the dynamicity that typifies many current distributed computing environments, e.g. grids. *muskel* allows the user to describe an application in terms of generic skeleton compositions. The description is then translated to a macro data flow graph [3] and the graph computed by a distributed data flow interpreter [2]. Central to the *muskel* system is the concept of a *manager* that is responsible for recruiting the computing resources used to implement the distributed data flow interpreter, distributing the fireable data flow instructions (tasks) and monitoring the activity of the computations.

While the performance results demonstrated the utility of *muskel*, it was noted in [2] that the centralized data flow instruction repository (taskpool) represented a bottleneck and the manager a potential single point of failure. The work reported on here addresses the latter of these issues. The planned reengineering of the *muskel* manager was seen as an opportunity to extend earlier

^{*}This research is carried out under the FP6 Network of Excellence CoreGRID funded by the European Commission (Contract IST-2002-004265).

related experiments [4] with the language Orc [5] to investigate if it could usefully be employed in the development of such management software. The intent was not to embark upon a full-blown formal development of a modified `muskel` manager (as was done earlier for the related *Lithium* system [6]), with attendant formulation and proof of its properties, but rather to discover what return might be obtained from the use of such a formal notation for *modest* effort. In this sense, the aim was in keeping with the lightweight approach to formal methods as advocated by, inter alia, Agerholm and Larsen [7].

Orc was viewed as being apt for two reasons. First, it is an orchestration language, and the job of the `muskel` manager is one of orchestrating computational resources and tasks; and, second, while there are many process calculi which may be used to describe and reason about distributed systems, the syntax of Orc was felt to be more appealing to the distributed system developer whose primary interest lies not in describing and proving formal properties of systems.

The approach taken was to reverse engineer the original `muskel` manager implementation to obtain an Orc description; attempt to derive, in semi-formal fashion, a specification of a modified manager based on decentralized management; and, use this derived specification as a basis for modifying the original code to obtain the decentralized management version of `muskel`. The work described in this paper is the first part of a more complex activity aimed at both removing the single point of failure represented by `muskel` manager *and* implementing a distributed data flow instruction repository, removing the current related bottleneck. While the second step is still ongoing, the first step provides a suitable vehicle to illustrate the proposed methodology.

2 `muskel`: an overview

`muskel` is a skeleton based parallel programming environment written in Java. The distinguishing feature of `muskel` with respect to other skeleton environments [8, 9] is the presence of an application manager. The `muskel` user instantiates a manager by providing the skeleton program to be computed, the input and the output streams containing the (independent) tasks to be computed and the results, respectively, and a performance contract modeling user performance expectations (currently, the only contract supported is the `ParDegree` one, requesting the manager to maintain a constant parallelism degree during application computation). The user then requests invocation of the `eval()` method of the manager and the application manager takes care of all the details relating to the parallel computation of the skeleton program.

When the user requires the computation of a skeleton program, the `muskel` system behaves as follows. The skeleton program is compiled to a macro data flow graph, i.e. a data flow graph of instructions modeled by significant portions of Java code corresponding to user `Sequential` skeletons [3]. A number of processing resources (sufficient to ensure the user performance contract) running an instance of the `muskel` run time are recruited from the network. The `muskel` run time on these remote resources provides an RMI object that can be used to

compute arbitrary macro data flow instructions, such as those derived from the skeleton program. For each task appearing on the input stream, a copy of the macro data flow graph is instantiated in a centralized `TaskPool`, with a fresh graph id [2]. A `ControlThread` is started for each of the `muskel` remote resources (`RemoteWorkers`) just discovered. The `ControlThread` repeatedly looks for a fireable instruction in the task pool (the data-flow implementation model ensures that all fireable instructions are independent and can be computed in parallel) and sends it to its associated `RemoteWorker`. That `RemoteWorker` computes the instruction and returns the results. The results are either stored in the appropriate data flow instruction(s) in the task pool or delivered to the output stream, depending on whether they are intermediate results or final ones. In the event of `RemoteWorker` failure, i.e. if either the remote node or the network connecting it to the local machine fails, the `ControlThread` informs the manager and it, in turn, requests the name of another machine running the `muskel` run time support from a centralized *discovery service* and forks a new `ControlThread` to manage it, while the `ControlThread` managing the failed remote node terminates after reinserting in the `TaskPool` the macro data flow instruction whose computation failed [1]. Note that the failures handled by the `muskel` manager are *fail-stop* failures, i.e. it is assumed that an unreachable remote worker will not simply restart working again, or, if it restarts, it does so in its initial state. `muskel` has already been demonstrated to be effective on both clusters and more widely distributed workstation networks and grids (see also www.di.unipi.it/~marcod/Muskel).

3 The Orc notation

The orchestration language Orc has been introduced by Misra and Cook [5]. Orc is targeted at the description of systems where the challenge lies in organising a set of computations, rather than in the computations themselves. Orc has, as primitive, the notion of a site call, which is intended to represent basic computations. A site, which represents the simplest form of Orc expression, either returns a *single* value or remains silent. Three operators (plus recursion) are provided for the orchestration of site calls:

1. operator $>$ (sequential composition)
 $E_1 > x > E_2(x)$ evaluates E_1 , receives a result x , calls E_2 with parameter x . If E_1 produces two results, say x and y , then E_2 is evaluated twice, once with argument x and once with argument y . The abbreviation $E_1 \gg E_2$ is used for $E_1 > x > E_2$ when evaluation of E_2 is independent of x .
2. operator $|$ (parallel composition)
 $(E_1 | E_2)$ evaluates E_1 and E_2 in parallel. Both evaluations may produce replies. Evaluation of the expression returns the merged output streams of E_1 and E_2 .
3. where (asymmetric parallel composition)
 E_1 where $x : \in E_2$ begins evaluation of both E_1 and $x : \in E_2$ in parallel.

Expression E_1 may name x in some of its site calls. Evaluation of E_1 may proceed until a dependency on x is encountered; evaluation is then delayed. The first value delivered by E_2 is returned in x ; evaluation of E_1 can proceed and the thread E_2 is halted.

Orc has a number of special sites:

- 0 never responds (0 can be used to terminate execution of threads);
- if b returns a signal if b is true and remains silent otherwise;
- $RTimer(t)$, always responds after t time units (can be used for time-outs);
- let always returns (publishes) its argument.

Finally, the notation $(i : 1 \leq i \leq 3 : worker_i)$ is used as an abbreviation for $(worker_1|worker_2|worker_3)$.

4 muskel manager: an Orc description

The Orc description presented focuses on the management component of `muskel`, and in particular on the discovery and recruitment of new remote workers in the event of remote worker failure. The compilation of the skeleton program to a data flow graph is not considered.

While Orc does not have an explicit concept of “process”, processes may be represented as expressions which, typically, name channels which are shared with other expressions. In Orc a channel is represented by a site [5]. $c.put(m)$ adds m to the end of the (FIFO) channel and publishes a signal. If the channel is non-empty $c.get$ publishes the value at the head and removes it; otherwise the caller of $c.get$ suspends until a value is available.

The activities of the processes of the `muskel` system are now described, followed by the Orc specification.

System The *system* comprises a program, pgm , to be executed (for simplicity a single program is considered: in reality a set of programs may be provided here); a set of tasks which are initially placed in a *taskpool*; a *discovery* mechanism which makes available processing engines (*remoteworkers*); and a *manager* which creates control threads and supplies them with remote workers. t is the time interval at which potential remote worker sites are polled; and, for simplicity, also the time allowed for a remote worker to perform its calculation before presumption of failure.

Discovery It is assumed that the call $g.can_execute(pgm)$ to a *remote worker* site returns its name, g , if it is capable of (in terms of resources) and willing to execute the program pgm , and remains silent otherwise. The call $rwork_erpool.add(g)$ adds the remote worker name g to the pool provided it is not already there. The *discovery* mechanism carries on indefinitely to cater for possible communication failure.

Manager The *manager* creates a number (*contract*) of control threads, supplies them with remote worker handles, monitors the control threads for failed

remote workers and, where necessary, supplies a control thread with a new remote worker.

Control thread A control thread (*ctrlthread*) repeatedly takes a task from the *taskpool* and uses its remote worker to execute the program *pgm* on this task. A result is added to the *resultpool*. A time-out indicates remote worker failure which causes the control thread to execute a call on an *alarm* channel while returning the unprocessed task to the *resultpool*. The replacement remote worker is delivered to the control thread via a channel, c_i .

Monitor The *monitor* awaits a call on the *alarm* channel and, when received, recruits and supplies the appropriate control thread, i , with a new remote worker via the channel, c_i .

$$\begin{aligned}
& \text{system}(pgm, tasks, contract, G, t) \triangleq \\
& \quad \text{taskpool.add}(tasks) \\
& \quad | \text{discovery}(G, pgm, t) \\
& \quad | \text{manager}(pgm, contract, t) \\
& \text{discovery}(G, pgm, t) \triangleq (|_{g \in G} (\text{if } remw \neq false \gg rworkerpool.add(remw) \\
& \quad \quad \quad \text{where } remw : \in \\
& \quad \quad \quad (\quad g.can_execute(pgms) \\
& \quad \quad \quad | \text{Rtimer}(t) \gg let(false)) \\
& \quad \quad \quad) \\
& \quad \quad \quad) \gg \text{discovery}(G, pgm, t) \\
& \text{manager}(pgm, contract, t) \triangleq \\
& \quad | i : 1 \leq i \leq contract : (rworkerpool.get > remw > ctrlthread_i(pgms, remw, t)) \\
& \quad | \text{monitor} \\
& \text{ctrlthread}_i(pgms, remw, t) \triangleq \text{taskpool.get} > tk > \\
& \quad (\quad \text{if } valid \gg \text{resultpool.add}(r) \gg \text{ctrlthread}_i(pgms, remw, t) \\
& \quad | \text{if } \neg valid \gg (\quad \text{taskpool.add}(tk) \\
& \quad \quad | \text{alarm.put}(i) \gg c_i.get > w > \text{ctrlthread}_i(pgms, w, t) \\
& \quad \quad) \\
& \quad) \\
& \quad \text{where } (valid, r) : \in \\
& \quad \quad (\quad remw(pgms, tk) > r > let(true, r) \\
& \quad \quad | \text{Rtimer}(t) \gg let(false, 0) \\
& \quad \quad) \\
& \text{monitor} \triangleq \text{alarm.get} > i > rworkerpool.get(remw) > remw > c_i.put(remw) \\
& \quad \quad \gg \text{monitor}
\end{aligned}$$

5 Decentralized management: derivation

In the *muskel* system described thus far, the *manager* is responsible for the recruitment and supply of (remote) workers to control threads, both initially and in the event of worker failure. Clearly, if the manager fails, then, depending on the time of failure, the fault recovery mechanism will cease or, at worst, the entire system of control thread recruitment will fail to initiate properly.

Thus, the aim is to devolve this management activity to the control threads themselves, making each responsible for its own worker recruitment.

The strategy adopted is to examine the execution of the system in terms of traces of the site calls made by the processes and highlight management related communications. The idea is to use these communications as a means of identifying where/how functionality may be dispersed. In detail, the strategy proceeds as follows:

1. Focus on communication actions concerned with management. Look for patterns based on the following observation. Typically communication occurs when a process, A, generates a value, x , and communicates it to B. Identify occurrences of this pattern and consider if generation of the item could be shifted to B and the communication removed, with the “receive” in B being replaced by the actions leading to x ’s generation. For example:

$$A : \dots a1, a2, a3, send(x), a4, a5, \dots$$

$$B : \dots b1, b2, b3, receive(i), b4, b5, \dots$$
 Assume that $a2, a3$ (which, in general, may not be contiguous) are responsible for generation of x , and it is reasonable to transfer this functionality to B. Then the above can be replaced by:

$$A : \dots a1, a4, a5, \dots$$

$$B : \dots b1, b2, b3, a2, a3, (b4, b5, \dots)_{[i/x]}$$
2. The following trace subsequences are identified:
 - In control thread: $\dots alarm.put(i) \gg c_i.get > e > ctrlthread_i(pgm, e, t) \dots$
 - In monitor: $\dots alarm.get > i > rworkerpool.get > e > c_i.put(e) \gg \dots$
3. The subsequence $rworkerpool.get > e > c_i.put(e)$ of *monitor* actions is responsible for generation of a value (a remote worker) and its forwarding to a *ctrlthread* process. In the *ctrlthread* process the corresponding “receive” is $c_i.get$. So, the two trace subsequences are modified to:
 - In control thread: $\dots alarm.put(i) \gg$
 $rworkerpool.get > e > ctrlthread_i(pgm, e, t) \dots$
 - In monitor: $\dots alarm.get > i > \dots$
4. The derived trace subsequences now include the communication of the control thread number, i from $ctrlthread_i$ to the *monitor*, but this is no longer required by *monitor*; so, this communication can be removed.
5. Thus the two trace subsequences become:
 - In control thread: $\dots \gg rworkerpool.get > e > ctrlthread_i(pgm, e, t) \dots$
 - In monitor: $\dots \gg \dots$
6. Now the specifications of the processes $ctrlthread_i$ and *monitor* are examined to see how their definition can be changed to achieve the above trace modification, and consideration is given as to whether such modification makes sense and achieves the overall goal.
 - (a) In monitor the entire body apart from the recursive call is eliminated thus prompting the removal of the *monitor* process entirely. This is

as would be expected: if management is successfully distributed then there is no need for centralised monitoring of control threads with respect to remote worker failure.

(b) In control thread the clause:

$$| \text{alarm.put}(i) \gg c_i.get > e > \text{ctrlthread}_i(\text{pgm}, e, t)$$

becomes

$$| \text{rworkerpool.get} > w > \text{ctrlthread}_i(\text{pgm}, w, t)$$

This now suggests that ctrlthread_i requires access to the rworkerpool . But the rworkerpool is an artefact of the (centralised) manager and the overall intent is to eliminate this manager. Thus, the action rworkerpool.get must be replaced by some action(s), local to ctrlthread_i , which has the effect of supplying a new remote worker. Since there is no longer a remote worker pool, on-the-fly recruitment of an remote worker is required. This can be achieved by using a discovery mechanism similar to that of the centralised manager and replacing rworkerpool.get by $\text{discover}(G, \text{pgm})$:

$$\text{discover}(G, \text{pgm}) \triangleq \text{let}(rw) \text{ where } rw : \in |_{g \in G} g.\text{can_execute}(\text{pgm})$$

(c) Finally, as there is no longer centralised recruitment of remote workers, the control thread processes are no longer instantiated with their initial remote worker but must recruit it themselves. This requires that

- i. the control thread process be further amended to allow initial recruitment of a remote worker, with the (formerly) recursive body of the process now defined within a subsidiary process, ctrlprocess , as shown below.
- ii. the parameter remw in ctrlthread be replaced by G as the control thread is no longer supplied with an (initial) remote worker, but must handle its own remote worker recruitment by reference to the grid, G .

The result of these modifications is shown in the decentralized manager specification below.

Decentralized Management Here each control thread is responsible for recruiting its own remote worker (using a discovery mechanism similar to that of

the centralised manager specification) and replacing it in the event of failure.

$$\begin{aligned}
& \text{systemD}(\text{pgm}, \text{tasks}, \text{contract}, G, t) \triangleq \\
& \quad \text{taskpool.add}(\text{tasks}) \\
& \quad | i : 1 \leq i \leq \text{contract} : \text{ctrlthread}_i(\text{pgm}, t, G) \\
& \quad \text{ctrlthread}_i(\text{pgm}, t, G) \triangleq \text{discover}(G, \text{pgm}) > rw > \text{ctrlprocess}(\text{pgm}, rw, t, G) \\
& \quad \text{discover}(G, \text{pgm}) \triangleq \text{let}(rw) \text{ where } rw : \in |_{g \in G} g.\text{can_execute}(\text{pgm}) \\
& \quad \text{ctrlprocess}(\text{pgm}, rw, t, G) \triangleq \text{taskpool.get} > tk > \\
& \quad \quad (\text{if } \text{valid} \gg \text{resultpool.add}(r) \gg \text{ctrlprocess}(\text{pgm}, rw, t, G) \\
& \quad \quad | \text{if } \neg \text{valid} \gg \text{taskpool.add}(tk) \\
& \quad \quad \quad | \text{discover}(G, \text{pgm}) > w > \\
& \quad \quad \quad \quad \text{ctrlprocess}(\text{pgm}, w, t, G) \\
& \quad \quad) \\
& \quad \text{where } (\text{valid}, r) : \in \\
& \quad \quad (\text{remw}(\text{pgm}, tk) > r > \text{let}(\text{true}, r) \\
& \quad \quad | \text{Rtimer}(t) \gg \text{let}(\text{false}, 0) \\
& \quad \quad)
\end{aligned}$$

5.1 Analysis

Having derived a decentralized manager specification, the “equivalence” of the two versions must be established. In this context, equivalent means that the same input/output relationship holds, as clearly the two systems are designed to exhibit different non-functional behaviour.

The input/output relationship (i.e. functional semantics) is driven almost entirely by the *taskpool*, whose contents change dynamically to represent the data-flow execution. This execution primarily consists in establishing an on-line partial order among the execution of fireable tasks. All execution traces compliant to this partial order exhibit the same functional semantics by definition of the underlying data-flow execution model. This can be formally proved by showing that all possible execution traces respecting data-dependencies among tasks are functionally confluent (see [6] for the full proof), even if they do not exhibit the same performance.

Informally, one can observe that a global order among the execution of tasks can not be established ex ante, since it depends on the program and the execution environment (e.g. task duration, remote workers’ availability and their relative speed, network connection speed, etc.). So, different runs of the centralized version will typically generate different orders of task execution. The separation of management issues from core functionality, which is a central plank of the *muskel* philosophy, allows the functional semantics of the centralized system to carry over intact to the decentralized version as this semantics is clearly independent of the means of recruiting remote workers.

One can also make an observation on how the overall performance of the system might be affected by these changes. In the centralised management system, the discovery activity is composed with the “real work” of the remote

workers by the parallel composition operator: discovery can unfold in parallel with computation. In the revised system, the discovery process is composed with core computation using the sequence operator, \gg . This suggests a possible price to pay for fault recovery.

6 Decentralized management: implementation

Following the derivation of the decentralized manager version outlined above, the existing `muskel` prototype was modified to introduce distributed fault management and to evaluate the relative cost in terms of performance. As shown above, in the decentralized manager, the $discovery(G, pgm, t)$ parallel component of the $system(\dots)$ expression become part (the $discover(G, pgm)$ expression) of the $ctrlprocess(\dots)$ expression. The $discovery$ and $discover$ definitions are not exactly the same, but $discover$ is easily derived from $discovery$. Thus, the code implementing $discovery(G, pgm, t)$ was moved and transformed appropriately to give an implementation of $discover(G, pgm)$. This required the modification of just one of the files in the `muskel` package (194 lines of code out of a total of 2575, less than 8%), the one implementing the control thread.

Experiments were run using the original and the modified versions to test the functionality and cost of the new implementation. The experiments were run on a Fast Ethernet network of Pentium III machines running Linux and Java 1.5. First the scalability of the decentralized manager version was verified. Figure 1 (upper plot) shows almost perfect scalability up to 8 nodes, comparable to that achieved when running the same program with the original `muskel`, both in the case of no faults and in the case of a single fault per computation. Then the times spent in managing a node fault in the centralized and decentralized versions were compared (Figure 1 lower part). The plot is relative to the time spent handling a single fault. The centralized version performs slightly better than the decentralized one, as anticipated. In the centralized version the discovery of the name of the remote machines hosting the `muskel` RTS is performed *concurrently* with the computation, whereas it is performed *serially* to the main computation in the decentralized version. The rest of the activities performed to handle the fault (lookup of the remote worker RMI object and delivery of the macro data flow) is the same in the two cases.

7 Conclusions

The manager component of the `muskel` system has been re-engineered to provide distributed remote worker discovery and fault recovery. A formal specification of the component, described in Orc, was developed. The specification provided the developer with a representation of the manager which allowed exploration of its properties and the development of what-if scenarios while hiding the inessential detail. By studying the communication patterns present within the process traces, the developers were able to derive a system exhibiting equivalent core

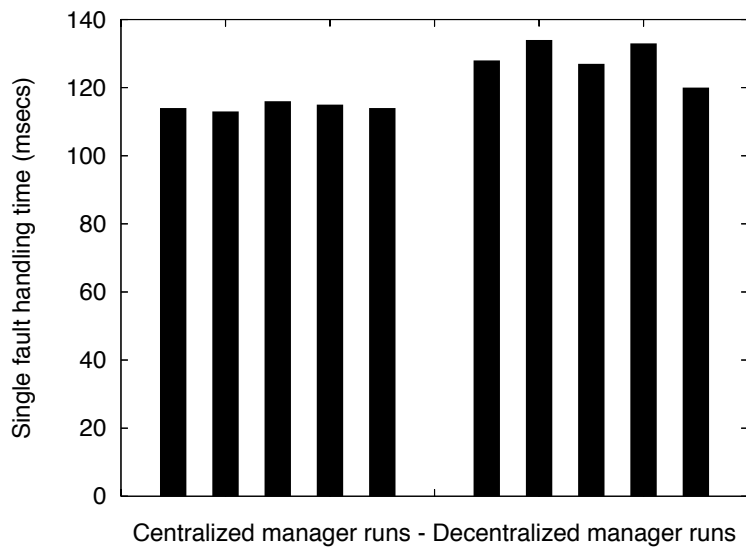
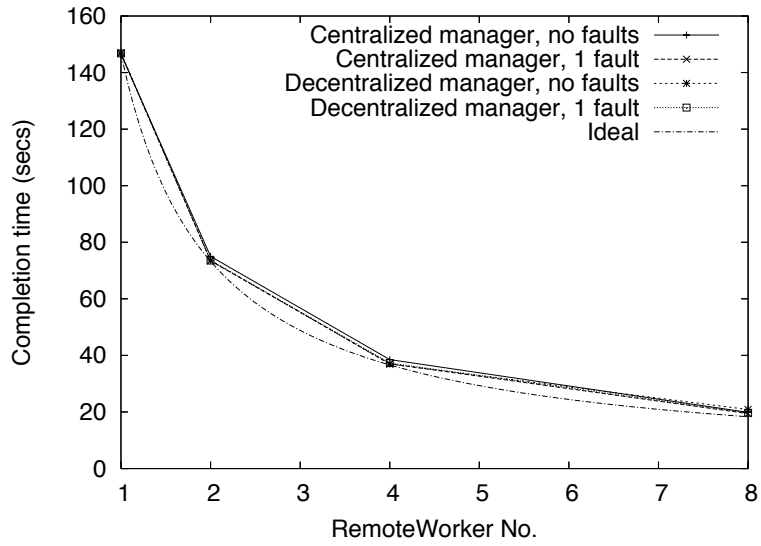


Figure 1: Scalability (upper) and fault handling cost (lower) of modified vs. original muskel

functionality, while having the desired decentralised management properties. The derivation proceeded in a series of semi-formally justified steps, with incorporation of insight and experience as exemplified by the inclusion of expressions such as “reasonable to transfer this functionality” and “such modification makes sense”.

The claim is that the creation of such a derivation facilitates exploration (and documentation) of ideas and delivers much return for small investment. Also, lightweight reasoning about the derived specification gave the developers some insight into the expected performance of the derived implementation relative to its parent.

Finally, the authors suggest that Orc is an appropriate vehicle for the description of management systems of the sort described here. Its syntax is small and readable; its constructs allow for easy description of the sorts of activities that typify these systems (in particular the asymmetric parallel composition operator facilitates easy expression of concepts such as time-out and parallel searching); and the *site* abstraction allows clear separation of management activity from core functionality.

Future work will involve tackling the more difficult task of removing the centralised task pool bottleneck, which should provide a stiffer test of the proposed approach. And, the availability of an Orc description makes possible the analysis of system variants with respect to cost and reliability using techniques described in [10].

References

- [1] Danelutto, M.: QoS in parallel programming through application managers. In: Proc. of Intl. Euromicro PDP: Parallel Distributed and network-based Processing, Lugano, Switzerland, IEEE (2005) 282–289
- [2] Danelutto, M., Dazzi, P.: Joint structured/non structured parallelism exploitation through data flow. In Alexandrov, V., van Albada, D., Sloat, P., Dongarra, J., eds.: Proc. of ICCS: Intl. Conference on Computational Science, Workshop on Practical Aspects of High-level Parallel Programming. LNCS, Reading, UK, Springer (2006)
- [3] Danelutto, M.: Dynamic run time support for skeletons. In D’Hollander, E.H., Joubert, G.R., Peters, F.J., Sips, H.J., eds.: Proc. of Intl. PARCO 99: Parallel Computing. Parallel Computing Fundamentals & Applications. Imperial College Press (1999) 460–467
- [4] Stewart, A., Gabarró, J., Clint, M., Harmer, T.J., Kilpatrick, P., Perrott, R.: Managing grid computations: An orc-based approach. In Guo, M., Yang, L.T., Martino, B.D., Zima, H.P., Dongarra, J., Tang, F., eds.: ISPA. Volume 4330 of LNCS., Springer (2006) 278–291

- [5] Misra, J., Cook, W.R.: Computation orchestration: A basis for a wide-area computing. *Software and Systems Modeling* (2006) DOI 10.1007/s10270-006-0012-1.
- [6] Aldinucci, M., Danelutto, M.: Skeleton based parallel programming: functional and parallel semantic in a single shot. *Computer Languages, Systems and Structures* (2006) DOI 10.1016/j.cl.2006.07.004, in press.
- [7] Agerholm, S., Larsen, P.G.: A lightweight approach to formal methods. In Hutter, D., Stephan, W., Traverso, P., Ullmann, M., eds.: *FM-Trends*. Volume 1641 of LNCS., Springer (1998) 168–183
- [8] Cole, M.: Bringing skeletons out of the closet: A pragmatic manifesto for skeletal parallel programming. *Parallel Computing* **30** (2004) 389–406
- [9] Kuchen, H.: The muesli home page (2006) <http://www.wi.uni-muenster.de/PI/forschung/Skeletons/>.
- [10] Stewart, A., Gabarró, J., Clint, M., Harmer, T.J., Kilpatrick, P., Perrott, R.: Estimating the reliability of web and grid orchestrations. In Gorlatch, S., Bubak, M., Priol, T., eds.: *Integrated Reserach in Grid Computing*, Kraków, Poland, CoreGRID, Academic Computer Centre CYFRONET AGH (2006) 141–152