



ELSEVIER

Information Sciences 140 (2002) 229–240

INFORMATION
SCIENCES

AN INTERNATIONAL JOURNAL

www.elsevier.com/locate/ins

Parallelising large irregular programs: an experience with Naira

Sahalu B. Junaidu ^{a,*}, Phil W. Trinder ^b

^a *Department of Information and Computer Science, King Fahd University of Petroleum and Minerals, Dhahran 31261, Saudi Arabia*

^b *Department of Computer Science and Electrical Engineering, Herriot-Watt University, Edinburgh, Scotland, UK*

Received 1 April 2000; received in revised form 8 May 2001; accepted 28 June 2001

Abstract

Naira is a compiler for Haskell, written in Glasgow parallel Haskell. It exhibits modest, but irregular, parallelism that is determined by properties of the program being compiled, e.g. the complexity of the types and of the pattern matching. We report four experiments into Naira's parallel behaviour using a set of realistic inputs: namely the 18 Haskell modules of Naira itself. The issues investigated are:

- Does increasing input size improve sequential efficiency and speedup?
- To what extent do high communications latencies reduce average parallelism and speedup?
- Does migrating running threads between processors improve average parallelism and speedup at all latencies? © 2002 Published by Elsevier Science Inc.

1. Introduction

The parallel behaviour of programs with regular parallelism is well understood: there are good simulation tools [10], a variety of analytical cost models [12,13] and hybrid methods [2]. Unfortunately many useful programs lack such a regular structure, for example the number and size of tasks may be deter-

* Corresponding author.

E-mail addresses: sahalu@ccse.kfupm.edu.sa (S.B. Junaidu), trinder@cee.hw.ac.uk (P.W. Trinder).

mined by the input. It is much harder to construct analytical models of the behaviour of programs with such irregular parallelism, and cost models are far less well developed [2,10,13]. Instead we must rely on simulation and measurement, and for these to be meaningful the program simulated must be realistic, and applied to real input data.

Parallel functional languages with dynamic models of parallelism are well suited to constructing programs with irregular parallelism. In the implementation of these languages many aspects of parallel behaviour, e.g. how many threads to create, and where to execute a thread, are automatically managed by a sophisticated runtime system. In consequence the runtime system can adapt at runtime to irregular parallel behaviour. Glasgow parallel Haskell (GpH) is one such language [11], and several irregularly parallel programs have been constructed [8].

This paper reports a series of experiments into the irregular parallel behaviour of the Naira compiler. Naira comprises approximately 5000 lines of GpH, and 1000 lines of C. As a compiler Naira performs symbolic manipulation of program text. Although the results of our experiments only apply directly to Naira, we hope that the principles established will hold for other symbolic programs with irregular parallelism. Naira is a suitable basis for these experiments because of the availability of a suite of real input data, and more importantly it is a large, real program that achieves wall-clock speedups, as reported in Section 2.

2. Naira

Naira compiles from a substantial subset of Haskell to a RISC-like target language that has been extended with special parallel constructs [9]. The front end of the compiler comprises about 5K lines of GpH code organised in 18 modules. Input is in the form of a Haskell module, i.e. a file containing function definitions. The phases of the compiler, and the data dependencies between them, are depicted in Fig. 1. The first phase, analysis, consists of the lexical analyser and the parser. The next four phases implement the pattern

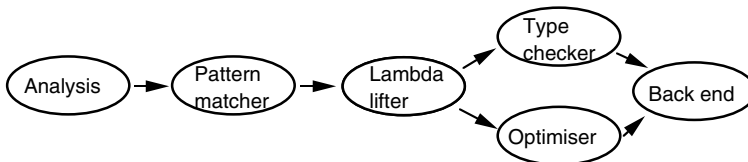


Fig. 1. Naira architecture.

matching compiler, the lambda lifter, the type checker and the intermediate language optimiser, respectively. The detailed organisation and implementation of these phases is described elsewhere [5].

The main sources of parallelism in Naira are as follows: data-parallel compilation of each function definition in the input; pipelining of data between the compilation phases identified in Fig. 1; control parallelism in the pattern matcher and lambda lifter.

Naira's modest parallelism has been measured on the GranSim simulator, and subsequently on a network of workstations. GranSim is a simulator that can be parameterised to emulate different parallel architectures [3]. The best simulated speedups achieved are 10.9 on an idealised parallel machine with an unbounded number of processors; 6.3 on an 8-processor shared-memory architecture; 5.8 on an 8-processor distributed-memory architecture. On a network of five Sun workstations Naira achieves a maximum speedup of 2.73 relative to its execution on a single workstation. However parallel GpH code carries an execution overhead compared to sequential code, and Naira's sequential efficiency is 89%. As a result the wall-clock speedup achieved 2.46 on five processors. These speedups agree with GranSim's predicted speedup of 3.01 for the workstation network. Detailed descriptions of the parallelisation and other aspects of Naira can be found in [4,5].

3. Experimental apparatus

The first experiment below measures Naira's performance using GpH's GUM runtime system. The other experiments use the GranSim parameterisable simulator [3] because it provides a cheap and convenient platform. GranSim is closely based on the GUM runtime system that manages execution of GpH programs on real parallel architectures. GranSim takes parameters describing many aspects of a parallel architecture, e.g. thread creation costs, communication latency, etc. As outlined above, the accuracy of GranSim predictions for Naira has been corroborated, and its accuracy for several other programs has been validated on several architectures [7].

The parallel program metrics used are standard [1], as follows:

Runtime: the time to execute a program.

Sequential efficiency: the ratio between runtime under the sequential runtime system, and runtime under the parallel runtime system on a single processor. Sequential efficiency measures the additional work that the parallel language must perform to synchronise with concurrent threads.

Average parallelism: the average number of active threads during a program's execution.

Speedup (at N processors): the ratio between sequential runtime and parallel runtime (on N processors).

4. Experiment 1: impact of input size

The first experiment investigates the impact of input size on both the sequential efficiency, and speedups achieved by Naira.

4.1. Input size vs sequential efficiency

What is the impact of input size on the sequential efficiency of symbolic irregularly parallel GpH programs? Table 1 shows the runtimes for Naira compiling each of its 18 GpH modules. The runtime figures in this table are the wall-clock (real) timings taken to compile each input module. These times (measured in seconds) are averaged over several program executions. The second column of Table 1 records the Naira runtime to compile of each module, when Naira itself is running, fully optimised, under the standard sequential runtime system. The third column gives the Naira runtime to compile

Table 1
Naira's runtime and sequential efficiency for each module

Input module	Sequential	Seq. for par. exec.	Efficiency (%)
MyPrelude	2.5	3.0	83
DataTypes	2.8	3.2	89
Tables	3.4	4.1	83
PrintUtils	1.7	2.5	70
Printer	3.3	4.0	82
LexUtils	3.4	3.8	89
Lexer	3.0	3.6	83
SyntaxUtils	22.2	23.3	95
Syntax	4.9	5.4	90
MatchUtils	4.5	5.2	86
Matcher	5.0	5.5	91
LambdaUtils	2.1	2.6	79
LambdaLift	7.2	8.2	88
TCheckUtils	10.0	11.8	84
TChecker	2.8	3.7	75
OptmiseUtils	7.7	9.4	82
Optimiser	10.6	12.6	84
Main	6.9	8.4	82
Minimum efficiency			70
Maximum efficiency			90
Mean efficiency			86

Table 2
Naira's runtime and sequential efficiency for 2-module inputs

Input module	Sequential	Seq. for par. exec.	Efficiency (%)
Preface	14.4	21.8	66
Printing	7.7	8.7	88
Lexing	8.4	9.1	92
Parsing	33.3	35.8	93
PMatching	13.2	14.6	90
LLifting	6.5	7.2	90
TChecking	15.5	17.9	86
Optimising	41.1	45.1	91
Minimum efficiency			66
Maximum efficiency			93
Mean efficiency			87

each module, under the parallel runtime system on a single processor. The fourth column contains the efficiency calculated by dividing the second column by the third, and shows that the overhead imposed by parallel execution can be high for some inputs, e.g. PrintUtils, or TChecker. The last three rows contain a summary of information in the table.

Naira's input size can be increased by concatenating modules, for example the module Preface in Table 2 is a result of concatenating the contents of the MyPrelude and DataTypes modules from Table 1. Table 2 reports the sequential efficiency of Naira compiling these larger inputs. Similarly Table 3 reports the sequential efficiency of Naira compiling files constructed by concatenating four of the original modules.

The effect of increasing input size can be observed by comparing Tables 1–3, and we note that the mean sequential efficiency of Naira modules increases from 86% to 91%. Maximum efficiency likewise increases uniformly, from 90% to 93%; however, it is not clear why minimal efficiency drops to 66% in Table 2. Sequential efficiency probably improves because fixed parallel exe-

Table 3
Naira's runtime and sequential efficiency for 4-module inputs

Input module	Sequential	Seq. for par. exec.	Efficiency (%)
MyPrelToPrinter	63.9	76.2	84
LexUToSyntax	67.1	74.8	90
MatchUToLLift	33.4	37.2	90
MAIN	115.3	120.3	96
Minimum efficiency			84
Maximum efficiency			96
Mean efficiency			91

cution overheads, like program start-up costs, are amortised over a longer runtime.

Conclusion 1.1. Sequential efficiency of Naira improves with input size.

4.2. Input size vs speedup

What is the impact of input size on the average parallelism and speedups achieved by symbolic irregularly parallel GpH programs? Naira's 18 GpH modules are a varied collection of real compiler inputs. Several size metrics are possible for compiler input, each with a different variability in the 18 modules. For example lines of source code, varying from 100 to 500 lines; number of function definitions, varying from 10 to 90 functions; number of output lines produced, varying from 100 to 1200 lines (measures input size in the sense that the number of output lines reflects the complexity of the functions compiled). The measure selected is the number of function definitions as it provides a suitably abstract view of compiler input.

Fig. 2 plots the graph for the size of 18 modules against speedup. We see a weak correlation between size and speedup. There are several reasons why speedup may be improved with larger input, most probably the overheads of parallel structures, like pipeline startup and shutdown, are amortised over a longer runtime. For irregularly parallel programs, however, a special factor may be at work: a longer runtime allows the runtime system more time to adapt, and take advantage of the adaptation.

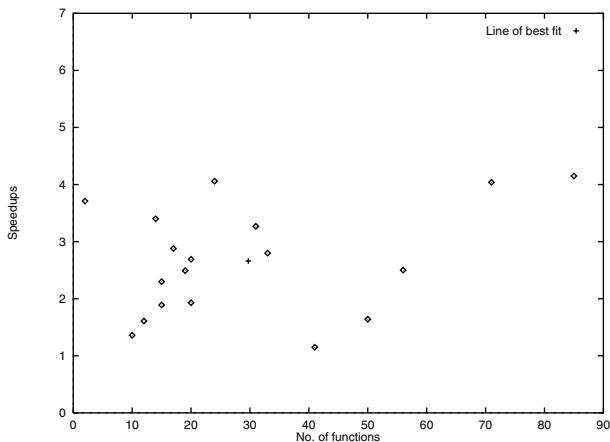


Fig. 2. Number of functions in input vs speedups.

Conclusion 1.2. There is a weak correlation between input size and speedup in Naira.

5. Experiment 2: impact of latency

To what extent do high communication latencies reduce average parallelism and speedup of symbolic irregularly parallel GpH programs? A key issue in parallel computing is achieving a high ratio of computation time to communication time. To discuss this issue in an architecture-independent way the time to send a message, or communications *latency*, is expressed in terms of the number of (computational) machine cycles. There are a whole range of parallel architectures, ranging from: shared-memory machines with a latency of around 5 or 10 cycles; fast distributed-memory machines (e.g. Meiko CS2) with a latency of 100–500 cycles; conventional distributed memory machines (e.g. IBM SP2) with a latency of 1000–5000 cycles; to networked machines (e.g. Suns on an ethernet) with a latency of 25 000–100 000 cycles.

Parallel language implementations in general, and GpH in particular, use a variety of mechanisms to hide latency costs, and high latencies expose the limitations of these mechanisms. GranSim allows us to measure a program on simulated architectures with varying latencies. Similar experiments have been reported elsewhere [3,6], but Naira is a larger program, and is measured on a set of real input data.

Table 4 and Figs. 3(a) and (b) summarise the average parallelism and speedups of Naira on a range of loosely coupled architectures with varying communication latencies. In the table the maximum and minimum average parallelism and speedup figures at each latency are boxed. We note that compiling some modules typically produces the greatest speedup, or the greatest average parallelism. Likewise some modules typically produce the least speedup, or average parallelism.

The table shows that the average parallelism and speedup for each program falls linearly as latency increases. Figs. 3(a) and (b) give a depiction of the trend for all programs, plotting mean, minimum and maximal speedup and average parallelism figures. This result is in contrast to earlier work on smaller irregular symbolic programs (approximately 500 lines) with greater speedups (approximately 28) that report an exponential reduction in speedup [3,6]. The difference may be because speedup and average parallelism are already low in Naira.

The figures also show that the difference between average parallelism and speedup is dramatic at a latency of 400 cycles; however as latency increases the difference between average parallelism and speedup diminishes. This indicates that at high latencies less unnecessary work is performed, but why this is so is currently unclear.

Table 4
Naira on 8-node GranSim with varying latencies: with migration

Input module	Communication latencies (distributed memory architectures)													
	400 cycles		2K cycles		5K cycles		25K cycles		50K cycles		85K cycles		120K cycles	
	Paral.	Spdup	Paral.	Spdup	Paral.	Spdup	Paral.	Spdup	Paral.	Spdup	Paral.	Spdup	Paral.	Spdup
MyPrelude	4.6	3.41	4.2	3.00	4.3	3.17	3.0	2.18	2.5	1.80	<u>1.8</u>	1.50	<u>1.6</u>	1.39
DataTypes	4.7	3.86	3.8	3.12	4.3	3.53	3.9	3.25	3.6	3.01	2.0	1.66	1.7	1.51
Tables	<u>5.0</u>	3.21	<u>4.6</u>	2.93	<u>4.9</u>	3.13	<u>4.1</u>	2.64	<u>4.0</u>	2.77	2.7	1.79	2.6	1.86
PrintUtils	<u>2.5</u>	2.66	<u>2.5</u>	2.62	<u>2.5</u>	2.57	<u>2.4</u>	2.52	<u>2.0</u>	2.09	1.9	2.06	1.8	2.04
Printer	<u>3.8</u>	3.98	3.8	3.90	3.7	3.87	3.3	<u>3.43</u>	<u>2.4</u>	2.52	2.6	2.71	1.8	1.88
LexUtils	<u>6.4</u>	3.18	5.6	4.32	5.4	3.32	5.1	<u>4.38</u>	<u>4.5</u>	3.89	3.7	2.99	2.7	2.50
Lexer	4.9	3.92	<u>4.7</u>	3.94	<u>4.5</u>	3.68	<u>3.2</u>	3.25	2.6	2.87	2.2	2.84	2.2	2.78
SyntaxUtils	6.2	2.46	<u>6.1</u>	2.54	<u>6.1</u>	2.31	<u>5.5</u>	2.10	3.7	4.33	3.6	4.55	2.7	3.19
Syntax	4.3	2.62	3.9	2.41	4.1	2.54	3.5	2.12	2.9	1.81	2.3	1.43	2.2	1.34
MatchUtils	3.9	4.05	3.6	3.75	3.7	3.86	2.9	3.04	3.1	4.06	1.8	1.87	1.7	2.04
Matcher	4.2	3.66	4.1	3.68	4.2	2.98	3.4	2.41	3.1	3.27	2.8	2.77	2.0	1.50
LambdaUtils	3.7	2.32	3.5	2.21	3.4	2.14	2.9	1.82	2.6	1.60	2.4	1.57	2.1	1.34
LambdaLift	4.6	<u>4.26</u>	4.3	<u>4.67</u>	4.0	<u>4.84</u>	3.4	3.06	2.8	<u>2.72</u>	<u>2.6</u>	<u>2.59</u>	1.9	<u>2.00</u>
TCheckUtils	5.1	<u>4.86</u>	5.2	<u>4.87</u>	5.2	<u>4.90</u>	4.3	<u>4.10</u>	4.3	<u>4.07</u>	<u>3.8</u>	<u>3.69</u>	<u>3.6</u>	<u>3.53</u>
TChecker	2.7	<u>1.77</u>	2.8	<u>1.71</u>	2.6	<u>1.77</u>	2.4	<u>1.48</u>	2.3	<u>1.49</u>	1.9	<u>1.17</u>	2.0	<u>1.27</u>
OptmiseUtils	3.0	3.59	3.0	3.55	2.6	3.11	2.4	2.92	2.3	2.72	2.0	2.33	1.8	2.23
Optimiser	3.4	4.76	3.3	4.67	3.3	4.59	3.0	4.23	2.9	4.05	2.3	3.88	2.0	3.66
Main	3.6	2.69	3.5	2.70	3.1	2.31	3.1	2.33	2.5	1.98	2.6	1.97	2.1	2.03
Mean	4.3	3.40	4.0	3.37	4.0	3.26	3.4	2.85	3.0	2.84	2.5	2.41	2.1	2.12

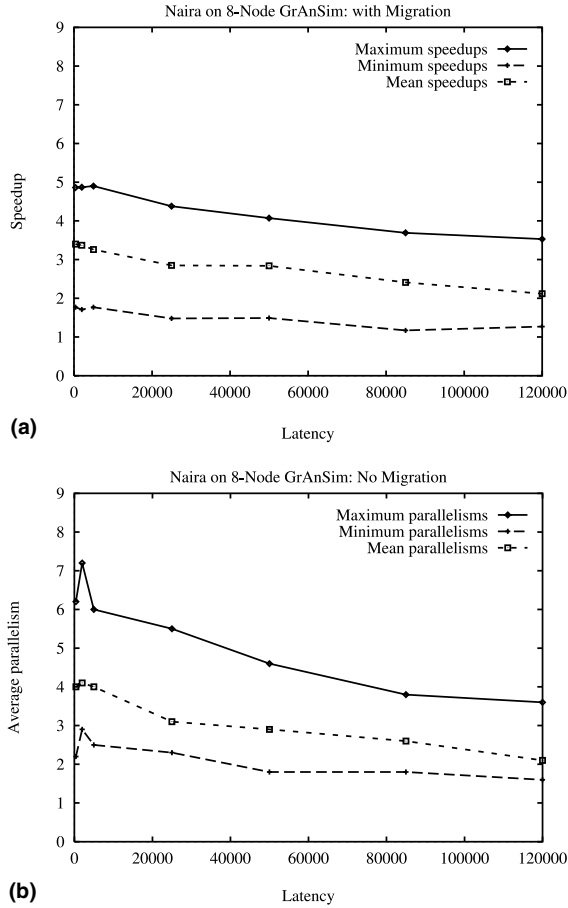


Fig. 3. Naira on 8-node GranSim: no migration. (a) Parallelism; (b) speedups.

Conclusion 2.1. Increased latency linearly reduces both average parallelism and speedup in Naira.

Conclusion 2.2. Increased latency reduces the difference between average parallelism and speedup in Naira.

6. Experiment 3: impact of thread migration

What is the impact of thread migration on the average parallelism and speedups achieved by symbolic irregularly parallel GpH programs at a range of latencies? Thread migration is an adaptive load-balancing technique for par-

allel languages where a thread of computation started on a heavily loaded processor may be moved to an idle processor. It is especially useful for correcting poor load distributions where one processor is overburdened with work. Implementing thread migration in a runtime is a complex undertaking, and it is not currently available in GpH. Part of the motivation for this experiment is to determine whether the benefits of thread migration for real programs are worth the implementation effort.

The program measurements reported in Table 4 have been repeated with GranSim parameterised to permit thread migration. The results are summarised in Figs. 4(a) and (b). At high latencies, i.e. above 5K cycles, both the speedups and average parallelism achieved with migration, and without it, are very similar. Differences are of the order of a few percent, typically in favour of with-migration, with occasional exceptions. This contradicts a body of earlier work on smaller irregular programs that indicates that thread migration gives sig-

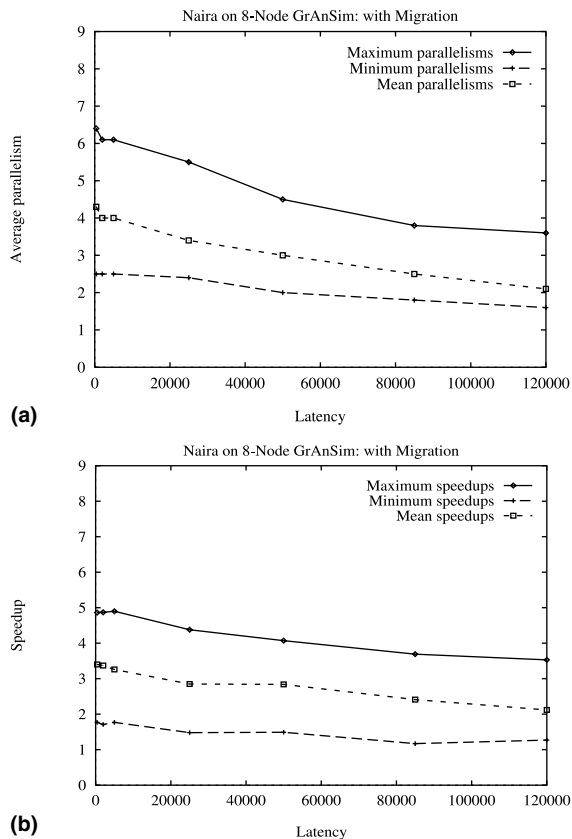


Fig. 4. Naira on 8-node GranSim: with migration. (a) Parallelism; (b) speedups.

nificant improvements [3]. There may be several reasons why migration does not improve Naira speedups at high latencies. Naira may be naturally generating a good load distribution, never requiring migration. Alternately the loss of data locality caused by migration may offset the benefits of migration, e.g. the cost of communicating large data structures like the symbol table outweigh the improved load distribution. Further investigation might establish whether Naira is typical of large irregular programs, or whether there are special considerations.

For most inputs at lower latencies speedup and average parallelism are similar with migration and without it. However, for a few inputs at lower latencies, speedup and average parallelism are dramatically improved by migration. This is reflected in the variability of Figs. 3(a) and (b) below 5K cycles, as compared with Figs. 4(a) and (b). Example modules with dramatic improvements are Syntax, and LexUtils, and we expect that poor load distributions are generated for these inputs. This result is consistent with the earlier work on thread migration [3], and extends it to medium-scale programs at low latencies. It also confirms earlier observations suggesting that for low-latency architectures a good load distribution is more important than good data locality [6].

Conclusion 3.1. At high latencies thread migration makes little difference to speedup or average parallelism in Naira.

Conclusion 3.2. At low latencies thread migration can make Naira's performance more predictable: giving dramatic improvements for some inputs.

7. Conclusions

We have measured several aspects of the irregularly parallel behaviour of the Naira compiler using the GranSim simulator. GranSim's prediction of Naira's behaviour has previously been validated on a network of workstations. We make the following conclusions: increasing input size improves both sequential efficiency and speedups. Raising communication latency reduces average parallelism and speedup linearly, and also reduces the difference between average parallelism and speedup. At low latencies we confirm earlier work on smaller programs showing that thread migration gives more predictable speedups and average parallelism, with dramatic improvements for some inputs. A new and more controversial result is that, for Naira at least, thread migration makes little difference at high latencies.

The experiments reported here are the first measurements of a medium-scale GpH symbolic and irregularly parallel program using a substantial set of real input data. We hope that future work will show that the principles established

above hold for a class of irregularly parallel programs. Further investigation may resolve some of the questions raised by the experiments, for example why does increased communication latency reduce the difference between average parallelism and speedup?

Acknowledgements

We would like to thank Tony Davie, Kevin Hammond, and Hans Wolfgang Loidl for stimulating discussions and other contributions to this work.

References

- [1] D.L. Eager, J. Zahorhan, E.D. Lazowska, Speedup versus efficiency in parallel systems, *IEEE Transactions on Computers* 38 (3) (1989) 408–423.
- [2] T. Fahringer, *Automatic Performance Prediction of Parallel Programs*, Kluwer Academic Publishers, Dordrecht, 1996.
- [3] K. Hammond, H.-W. Loidl, A.S. Partridge, Visualising granularity in parallel programs: a graphical winnowing system for Haskell, in: *HPFC'95 – Conference on High Performance Functional Computing*, Denver, CO, April 10–12, 1995, pp. 208–221.
- [4] S. Junaidu, A. Davie, K. Hammond, Naira: a parallel² Haskell compiler, in: *9th International Workshop on Implementing Functional Languages*, London, UK, September 1997, pp. 214–230.
- [5] S. Junaidu, *A parallel functional language compiler for message passing multicomputers*, Ph.D. Thesis, School of Mathematical and Computational Sciences, St. Andrews University, Scotland, 1998.
- [6] H.-W. Loidl, K. Hammond, Making a packet: cost-effective communication for a parallel graph reducer, in: *IFL'96 – International Workshop on the Implementation of Functional Languages*, Bad Godesberg, Germany, September 1996, *Lecture Notes in Computer Science*, vol. 1268, Springer, pp. 184–199.
- [7] H.-W. Loidl, *Granularity in large-scale parallel functional programming*, Ph.D. Thesis, Department of Computer Science, University of Glasgow, 1998.
- [8] H.-W. Loidl, P.W. Trinder, K. Hammond, S.B. Junaidu, R.G. Morgan, S.L. Peyton Jones, Engineering parallel symbolic programs in GpH, *Concurrency Practice and Experience* 11 (12) (December 1999) pp. 701–752.
- [9] G. Ostheimer, *Parallel functional programming for message passing multiprocessors*, Ph.D. Thesis, Department of Mathematical and Computational Sciences, St. Andrews University, Scotland, 1993.
- [10] V. Sarkar, *Partitioning and Scheduling Parallel Programs for Multiprocessors*, MIT Press, Cambridge, MA, 1989.
- [11] P.W. Trinder, K. Hammond, J.S. Mattson Jr., A.S. Partridge, S.L. Peyton Jones, GUM: a portable parallel implementation of Haskell, in: *PLDI '96 – Programming Languages Design and Implementation*, Philadelphia, PA, May 1996, pp. 78–88.
- [12] K.Y. Wang, A framework for static, precise performance prediction for superscalar-based parallel computers, in: *4th International Workshop on Compilers for Parallel Computers*, Delft, The Netherlands, December 1993, pp. 413–427.
- [13] N. Yazici-Pekergin, J.M. Vincent, Stochastic bounds on execution times of parallel programs, *IEEE Transactions on Software Engineering* 17 (10) (1991) 1059–1068.