

Iterative compilation in a non-linear optimisation space

François Bodin, Toru Kisuki, Peter Knijnenburg, Mike O' Boyle, Erven Rohou

► **To cite this version:**

François Bodin, Toru Kisuki, Peter Knijnenburg, Mike O' Boyle, Erven Rohou. Iterative compilation in a non-linear optimisation space. Workshop on Profile and Feedback-Directed Compilation, Oct 1998, Paris, France. 1998. <inria-00475919>

HAL Id: inria-00475919

<https://hal.inria.fr/inria-00475919>

Submitted on 23 Apr 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Iterative Compilation in a Non-Linear Optimisation Space

F. Bodin * T. Kisuki † P.M.W. Knijnenburg † M.F.P. O’Boyle ‡ E.Rohou*

Abstract

This paper investigates the applicability of iterative search techniques in program optimisation. Iterative compilation is usually considered too expensive for general purpose computing but is applicable to embedded applications where the cost is easily amortised over the number of embedded systems produced. This paper presents a case study, where an iterative search algorithm is used to investigate a non-linear transformation space and find the fastest execution time within a fixed number of evaluations. By using profile feedback in the form of execution time, it searches a large but restricted transformation space and shows performance improvement over existing approaches. We show that in the case of large transformation spaces, we can achieve within 0.3% of the best possible time by visiting less than 0.25% of the space using a simple algorithm and find the minimum after visiting up to less than 1 % of the space.

1. Introduction

The use of transformations to improve program performance has been extensively studied for over 30 years. Such work is based primarily on static analysis, possibly with some profile information to determine the significant regions of the code [7] and runtime dependent control-flow. Each technique is characterised by trying (i) to determine how such a program would perform on a particular processor and (ii) developing a program transformation such that the code is more likely to execute more efficiently. Such an approach relies on modeling those features of the architecture and the program that are considered important. Although there has been improvement, much work remains because it is extremely difficult to accurately model program/machine interaction. The problem is made worse in that the number of transformations to apply is potentially infinite and therefore the transformation space considered must be highly restricted.

This paper examines another approach to this problem by first describing the problem as that of searching a non-linear optimisation space in order to find the minima. Analytic techniques typically consider one parameter in the optimisation space at a time i.e. tile size, data layout etc. to avoid dealing with combinatorial problems. In practise, however we show that program transformations are not independent in their effect on performance but are in fact closely coupled. Furthermore, the optimisation space is highly non-linear (see figure 1) with many local minima. In order to search this space we propose a technique based on iterative compilation. Different transformations are applied, corresponding to points in the transformation space, and their worth evaluated by executing the program. In order that such an approach is feasible, we need to minimise the number of points evaluated. We show that in the case of large transformation spaces, we can achieve within 0.3% of the best possible time by visiting less than 0.25% of the space using a simple algorithm and find the minimum after visiting less than 1 % of the space.

Although such an approach is usually ruled out in terms of excessive compilation time ¹, it is applicable in those instances where the same application is to be executed many times. Embedded systems are an extreme example of this and in the following section we briefly describe a compiler framework developed to optimise multimedia codes for embedded systems. In fact there exists a spectrum of techniques from those based purely on static analysis, to the approach described here which makes few assumptions about the program/processor. It is our intention to integrate these approaches by using the results of different static analysis as “seed points” in the optimisation space i.e. they will form part of the initial set of points with which to start the iterative evaluation. This guarantees that the iterative approach will find at least as good a solution as static approaches.

This paper describes a case study to see if iterative compilation can be worthwhile. A simple, well studied problem, matrix multiplication is selected and executed on four different processors. The impact of transformations on performance is described by enumerating all the points in the transformation space considered. This is followed by the

* IRISA, Campus Universitaire de Beaulieu, 35042 Rennes, France

† Dept, of Computer Science, Leiden University, The Netherlands

‡ Division of Informatics, The University, Edinburgh

¹compilation time now includes several runs of the program

description of a simple search algorithm that tries to find the best performance within the fewest number of evaluations. This is repeated for a larger transformation space and the behaviour of the search algorithm evaluated. The algorithm is then finally applied to the TriMedia-1000 simulator (a VLIW processor produced by Philips aimed at the embedded processor market) and its behaviour evaluated. The paper finishes with a brief survey of related work and some concluding remarks.

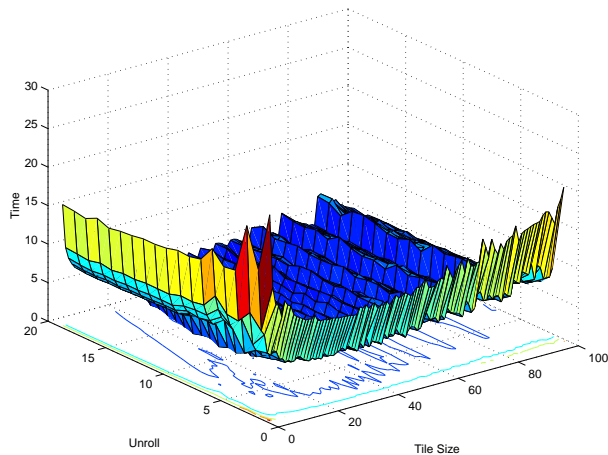


Figure 1. Transformation Space: UltraSparc N = 512

2. Iterative Compilation

The reason that program optimisation is difficult is that we are trying to minimise a function (execution time) which is undecidable at compile time over the infinite space of transformations. One solution is to actually evaluate the program at certain points in this space iteratively. Iterative compilation is generally not considered a viable proposition due the size of the space to be searched and hence the excessive compilation time. This is not the case for embedded systems where there is just one application to be optimised and long compilation times are acceptable if they increase performance. As a guide executing MPEG-2 for a reasonable number of frames takes less than 30 seconds on a TriMedia-1000. It typically takes an application programmer 3 months to produce an efficient implementation in which time a quarter of a million versions could be evaluated - more if there is more than one platform.

Analytic techniques have produced good results but are limited to the number of parameters they can consider. The size of the transformation space must be limited and assumptions about the low-level compiler made. Furthermore if a part of the system changes that is not implicitly or

explicit modeled - such as the register allocation policy of the local compiler, then analytic approaches are unable to adapt. Other factors generally ignored, include the resulting program structure when applying transformations with parameters relatively prime to the loop bounds, the introduction of spill code and instruction cache misses. Iterative compilation by definition consider all parts of the system when deciding on the best optimisation.

The Oceans project is an ESPRIT funded project, concerned with developing an iterative compiler for embedded systems. In particular we are targeting general purpose VLIW processors of which the Philips TriMedia-1 is typical. Economies of scale allow the production of cheaper and faster general purpose processors over custom embedded processors. However, such processors rely on efficient software implementations of the embedded applications but can afford long compilation times, hence our interest in iterative compilation.

A compiler infrastructure of high-level restructurer, code generator and low-level scheduler has been developed [1] for such processors and current work is focusing on the use of feedback/profile information in determining program optimisation. It is the long term goal of this project, to successfully integrate static analysis and feedback information for embedded application performance.

3. Problem Description

This paper is concerned with the investigation of iterative compilation for program performance. This is to be performed by way of a case study. Matrix multiplication is selected as the program to optimise because (i) it is well known - allowing independent comparison with other techniques (ii) it forms the core of the fdct in MPEG-2 and (iii) there are several legal transformations that can be applied. Furthermore, if we can show improvement for this extremely well-studied problem, then it is likely further improvement will be possible in those programs having received less attention. Other programs are currently being investigated, and will be presented in the final version of this paper.

The parameters of our experiment are as follows:

- Problem Size: $N = 400, 512$
- Processor Type: UltraSparc, R10000, Pentium Pro, and TriMedia -1000
- Transformation Space : Loop Unrolling 1-20, Tiling 1-100, Padding 1-10

We are interested in the correct combination of optimisations that minimises execution time on each processor for each data size. In this work we only consider high level

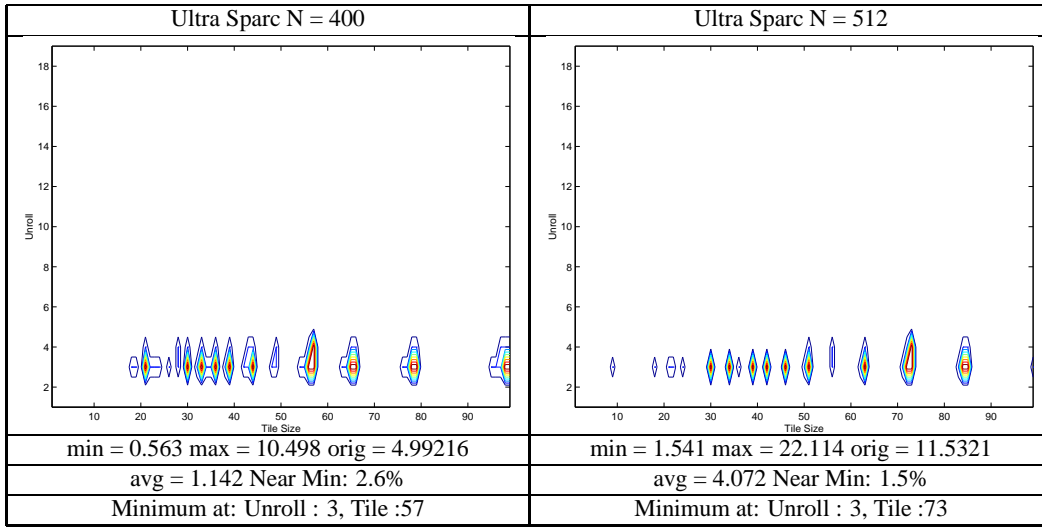


Figure 2. Transformation Space Characteristics of UltraSparc

optimisations and the only feedback information used is execution time. Further work will also include investigation of low-level optimisations and information regarding slot utilisation, register pressure etc.

In the following sections, the properties of the optimisation space are described and the impact of our algorithm in finding minima. Initially, we consider just tiling and unrolling for existing commodity processors. We then extend the experiments to consider padding and then apply the algorithm to the TriMedia-1000 simulator.

4. Tiling + Unrolling

Figure 1 shows a typical transformation space; in this case the performance of the UltraSparc for N=512 and varying tile and unroll size. It is periodic with high frequency oscillation and many local minima. Within such spaces, it is difficult to find the absolute minimum. The graphs in figures 2 to 5 show the percentage of the space within 20% of the minimum and its distribution, for the four different commodity processors and two different data sizes. The original², minimum, maximum and average execution for the space is also shown as is the number of points near 20% of the minimum execution time. What is immediately apparent is that the best transformation depends largely on the processor and to a lesser extent on the data size. The UltraSparc and Alpha perform best with a small unroll factor but the tile size varies. The Pentium has a more dispersed range of minima while the R10000 has the largest percentage of points closest to the minimum. Such a characteristic should

²All programs including transformed ones were compiled with -O2 optimisation

increase the probability of finding a good result within a reasonable number of samples. By way of contrast, the number of minimal points is much smaller for the Alpha and there is a much wider range of performance values - and incorrect optimisation may cost dearly. Across the various examples, optimisation gives an improvement of between a factor of 1.8 and 10 over unoptimised code.

5. Iterative Strategy

As can be seen from the the graphs in figures 2 to 5, there are several local minima per application and the minima varies from one data size/processor to the next. Clearly any algorithm that wishes to search the space must be robust enough so as not to be trapped in a local minimum. Therefore, techniques based on gradient approaches are not applicable. This must be balanced against searching too much of the transformation space, especially those regions where no suitable candidates can be found. This problem is compounded by the occurrence of minimal points surrounded by large values. It is not the purpose of this paper to propose the best search algorithm for such spaces, instead we are interested in the applicability of such algorithms.

Our search algorithm visits a number of points at spaced intervals, applying the appropriate transformation, executing the transformed program and evaluating its worth by measuring the execution time. Those points lying between the current global minimum and the average are added to an ordered queue. Iteratively, such points are removed from the queue and points within the neighbouring region are investigated, again at spaced intervals. This process is continued until a specific number of points have been evaluated

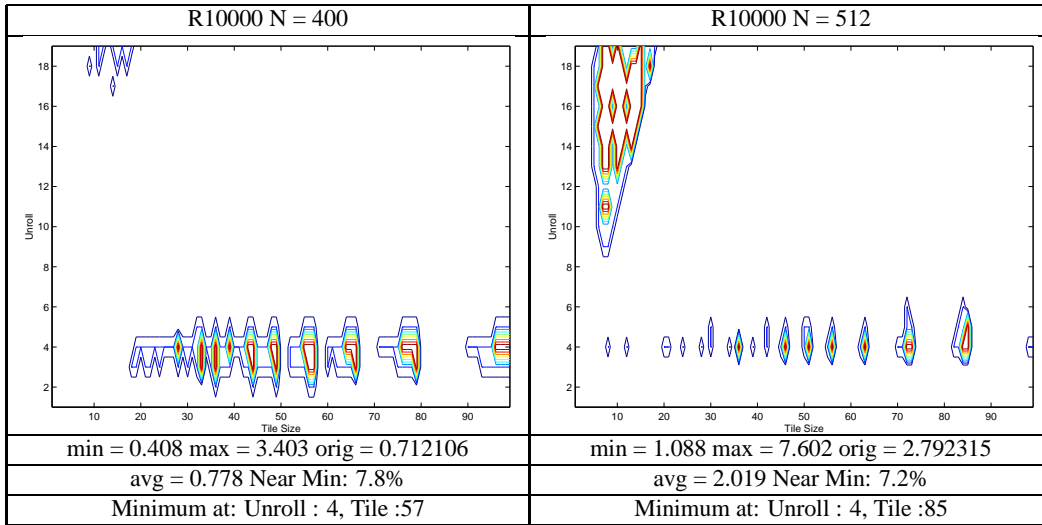


Figure 3. Transformation Space Characteristics of R10000

whereupon the point with minimal value -i.e. the fastest transformed program is reported.

5.1. Step Size vs Iteration Count

The step size within each of the transformation dimensions is a key component of the search algorithm. The step size that gives the biggest improvement depends not only on the program/data size/processor but also on the total number of evaluations to be undertaken. For instance a particular step size may be best if we are considering just 20 evaluations, but if we increase this to 200, another step size may be preferable. Despite this relationship, it was found that an initial of five - i.e. five samples in any one dimension gave a reasonable performance regardless of the number of evaluations undertaken or processor/data size/transformations considered.

5.2. Performance

Figure 6 and 7 gives the performance results of 4 of the 8 examples. The remaining 4 have similar performance. The x-axis of each graph is the number of evaluations carried out for each application of the iterative algorithm. The y-axis shows the percentage difference from the absolute minimum, the current global minimum is. For instance in the first graph after just one evaluation, the best performance found is 165% more than the actual minimum. i.e it is 2.65 times slower than the actual minimum. In the case of the Alpha it is 5.25 times slower. In each case the search algorithm finds the minimum in less than 200 steps or 10 % of the search space. More importantly, within 20 steps or 1% of the space, it is within 23.75 % of the minimum . There is

a rapid improvement in performance up to about 20 steps, with a more gradual improvement later. Despite the difficult characteristics of the transformation space, a relatively straightforward optimisation algorithm can find a very good implementation with a relatively few number of evaluations.

6. Padding

In order to further investigate the worth of iterative compilation, we extend the search space to include array padding in the first dimension of each array up to a pad size of 10, increasing the possible number of transformations to consider to 20 000. Clearly, in practice, it is unrealistic to evaluate exhaustively the transformation space, but in order to evaluate the performance of an iterative algorithm, we generated all possible program transformations for just two processors, the UltraSparc and Alpha for the data size N=512. The evaluation space is 4 dimensional and cannot be easily be presented graphically. Instead, certain slices of the space are shown in figures 8 and 9. Without padding, the best transformations, for the UltraSparc, are for the case when unroll = 3. When padding is also considered, however, none of the points within 20 % of the minima have an unroll factor of 3. This demonstrates the close connection of transformations and the error introduced when considering them separately. In fact the majority of minimal points occur when unrolling 4 times as shown in figure 8. Without padding, the best tile size was found to be 73 and in figure 8 the near minimal points for this tile size are shown. In fact, the best performance found is when the tile size is 51, unrolling 4 times and a pad size of 8. The execution time is now 1.2269 - a 20 % improvement in execution time when

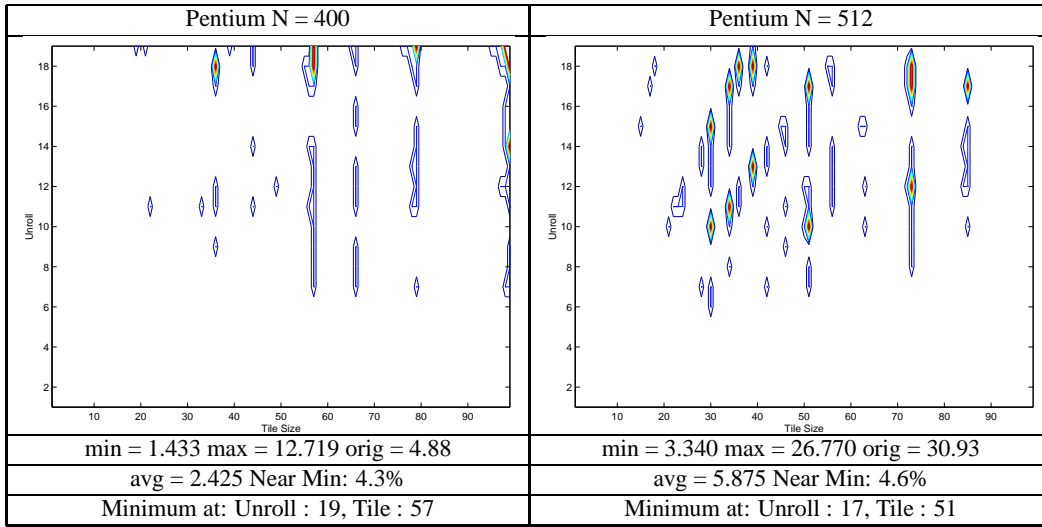


Figure 4. Transformation Space Characteristics of Pentium

compared to just unrolling and tiling. Solving one or more transformations independently will not give the best overall combination of transformations. In the case of the Alpha processor, shown in figure 9, unrolling by a factor of 4 still gives the best improvement, but now with a tile size of 73 and padding of size 8. As in the case of the UltraSparc, the improvement is an approximate 20 % reduction in execution time.

Although the number of points in the space has increased by a factor of 10 to 20000, the number of points needed to be evaluated has not grown proportionally. In fact within 50 evaluations we have found a transformed program while performance is within 0.3 % of the global minima on the UltraSparc. This can be seen in figure 10 where there is a rapid improvement in performance from over 7 times the minimal execution time to 0.3 % of the minima in 42 steps. Also shown in figure 10, is the same plot, with the first 20 iterations removed so as to give more detail on the algorithms performance. A similar behaviour is also seen in figure 11 for the Alpha. It, however, needs just over 82 steps to approach within 2% of the minimal possible execution time, finding the actual minima within 136 evaluations.

7. TriMedia-1000

The previous sections have described the transformation space and the performance of a search algorithm for several commodity processors. We are particularly interested in applying such techniques to embedded processors such as the TriMedia-1000 and in this section we evaluate our iterative approach to compilation by running the selected transformed programs on a cycle accurate simulator. As

we are using a simulator rather than actual processor, cycle counts are given as the performance measure and smaller data sizes are considered, namely $N=64$ and $N=128$.

7.1. Performance

Figure 12 shows those points within 20 % of the found minimum and their distribution for the case $N=64$. The performance of the search algorithm is also given. What is immediately noticeable, is that the number of points needed to sample does not scale down linearly with the size of the transformation space, just as it did not scale up when considering padding (see section 6.).

Although this paper has evaluated all points in the transformation space to give an absolute measure on the performance improvement of the iterative algorithm, in practice no absolute minima will be available. Rather the scheme will evaluate points returning the best available as long as sufficient time remains. This is the case with the final experiment where $N = 128$ on the TriMedia-1000. Due to the excessive simulation time it is not feasible to exhaustively search the space, so no absolute measure of performance is available. Nevertheless, the results in figures 13 show that the iterative algorithm makes steady improvement, reducing the execution time by a factor of 7 over the original program in less than 60 evaluations.

8. Related Work + Discussion

There is a large body of working considering program transformations to improve uniprocessor performance. In [2], an analytic algorithm to give a good tile size to min-

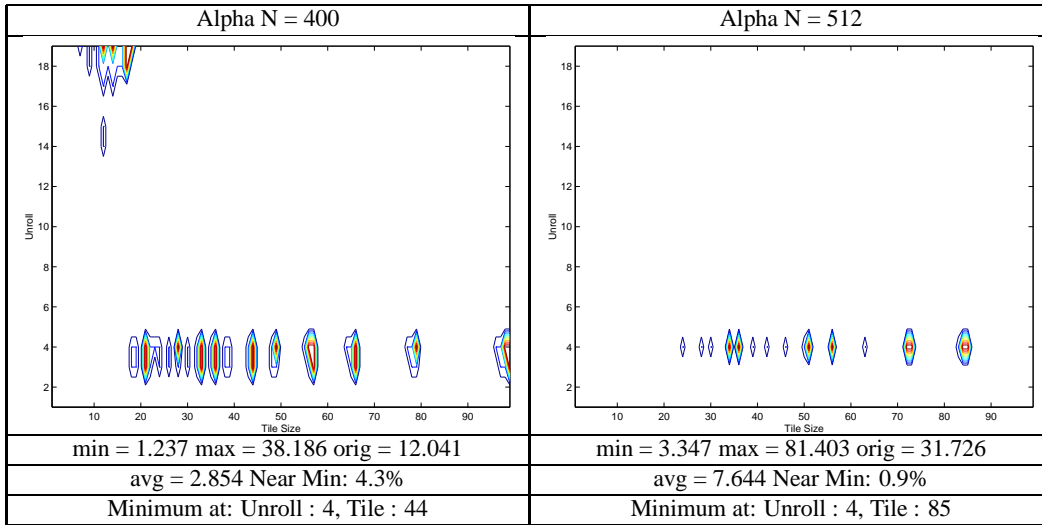


Figure 5. Transformation Space Characteristics of Alpha

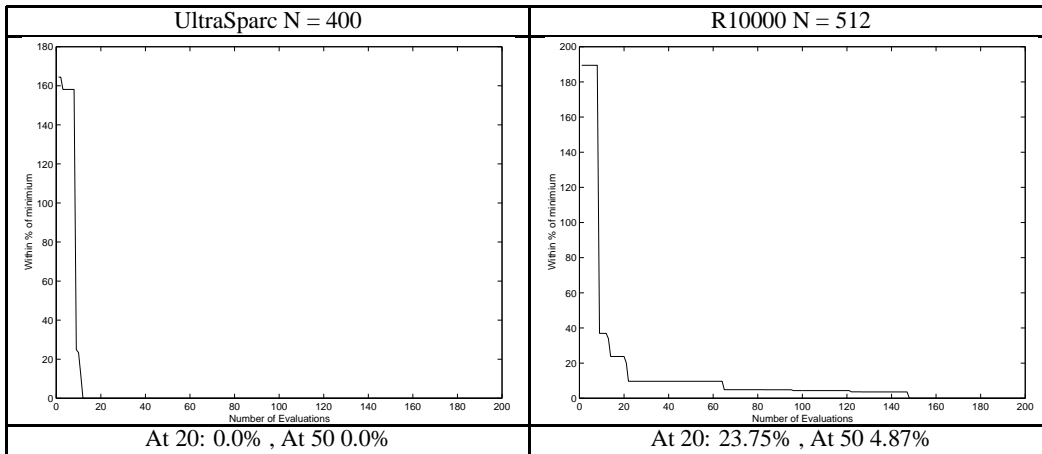


Figure 6. Search Algorithm Performance

imise interference and exploit locality is presented. This work considered rectangular tiles whose dimensions are a function of the iteration space and the cache organisation. This work gives good performance improvements over existing techniques but is wholly static and does not consider the impact of tiling on unrolling or other transformations. For example in [2], a tile size of 170×2 was considered to give the best performance on the Alpha for matrix multiplication when $N = 256$. Applying such a transformation gives an execution time of 0.794460 seconds. On applying our iterative algorithm, however, we reduce this time to 0.33863 with a tile size of 17, unrolling 18 times and array padding of 8. Static analysis could be used to “seed” the transformation space - i.e. give initial points with which to investigate. As they generally consider just one or two

transformations they will form search hyper planes with potentially many points to investigate. Further analysis techniques could be used to reduce their number. In this paper we have considered execution time as the metric for evaluating goodness of a transformation. As our system [1] provides additional information such as code size, register pressure, slot utilisation etc., it is possible to statically evaluate the goodness of a transformation after code generation. Although only approximate, as cache effects etc., cannot be exactly determined, such information may be used to prune transformed programs guaranteed to perform poorly.

Several researchers have considered using runtime information to select the best implementation. They however define one or more options statically which are then considered at runtime. For example in [4], whether or not a

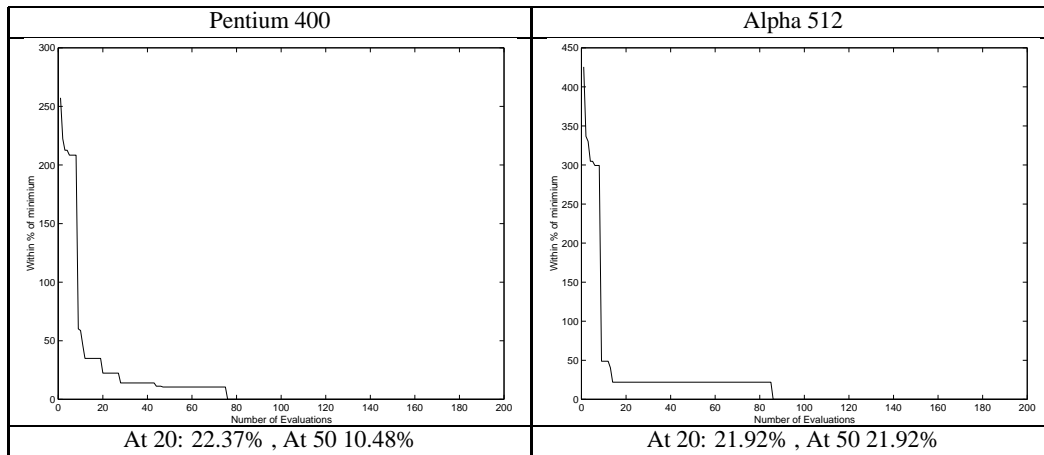


Figure 7. Search Algorithm Performance

portion of the iteration space should be tiled depends on runtime characteristics and in [6], different synchronisation algorithms are called depending on runtime behaviour. The work in this paper however, considers a much larger space of optimisations at compile time without incurring runtime overhead. Later work could combine the approaches by including dynamic monitoring to select, at runtime, one of a number of optimisations programs that were determined (at compile time) to perform well under certain circumstances - using the search space approach described in this paper. In [5] genetic algorithms are used to create and select transformations for parallel optimisation. This work is similar in spirit to the work presented in this paper but at present generates many illegal programs which must be discarded and hence examines a much larger set of programs before finding any improved solutions.

In this paper we have used a very simple search algorithm as a basis for iterative compilation. There is in fact a large literature on non-linear optimisation [3] though it is based on a continuous underlying optimisation function rather than the discrete space we consider. Techniques such as polynomial fitting could be applied to help improve the performance of the search algorithm. Although the best transformations to select are interdependent, knowledge of the processor and application domain could help bias the search space to first consider the areas where the minima is most likely to exist. For instance, in several of the examples, unrolling near a factor of 4 often leads to good results - however this is not always the case e.g. Pentium (see figure 4).

Although the search spaces considered are large, the domain is extremely limited. Only 3 transformations were applied to a simple loop nest. For general programs, the space to consider will be much greater. Future work needs to consider the application of an iterative compilation to much

larger programs. This paper has also focussed on parameterised transformations - but many transformations are not parameterised, they are either applied or not. For such transformations, a multi-dimensional search space may not be an appropriate representation and maybe search trees are more appropriate.

While it may be feasible to efficiently search large transformation spaces it is not necessarily always easy to generate the transformed program. Applying a sequence of program transformations corresponding to the position in the space is limited by the form of code produced by any previous transformation. In future, it will therefore be necessary to develop program transformation techniques to allow general iterative compilation consisting of many compound transformations.

This paper has concentrated on the effect of temporal performance. However, in embedded systems, code size is also important as it determines the amount of ROM required [8]. If we have a cost metric which is a function of execution time and code size, the techniques described in this paper can be immediately applied.

Finally, although we have focussed on maximising performance for a particular processor, the processor concerned strongly affects the likelihood of finding a good solution. Those processors which have less non-linear optimisation spaces are likely to be easier to optimise for.

9. Conclusion

This paper has investigated the use of iterative compilation as a program optimisation technique. By way of a simple case study, it has shown that determining the best transformations statically for a processor/data size pair is non-trivially. Furthermore, by considering program optimisation as searching a transformation space for minima it

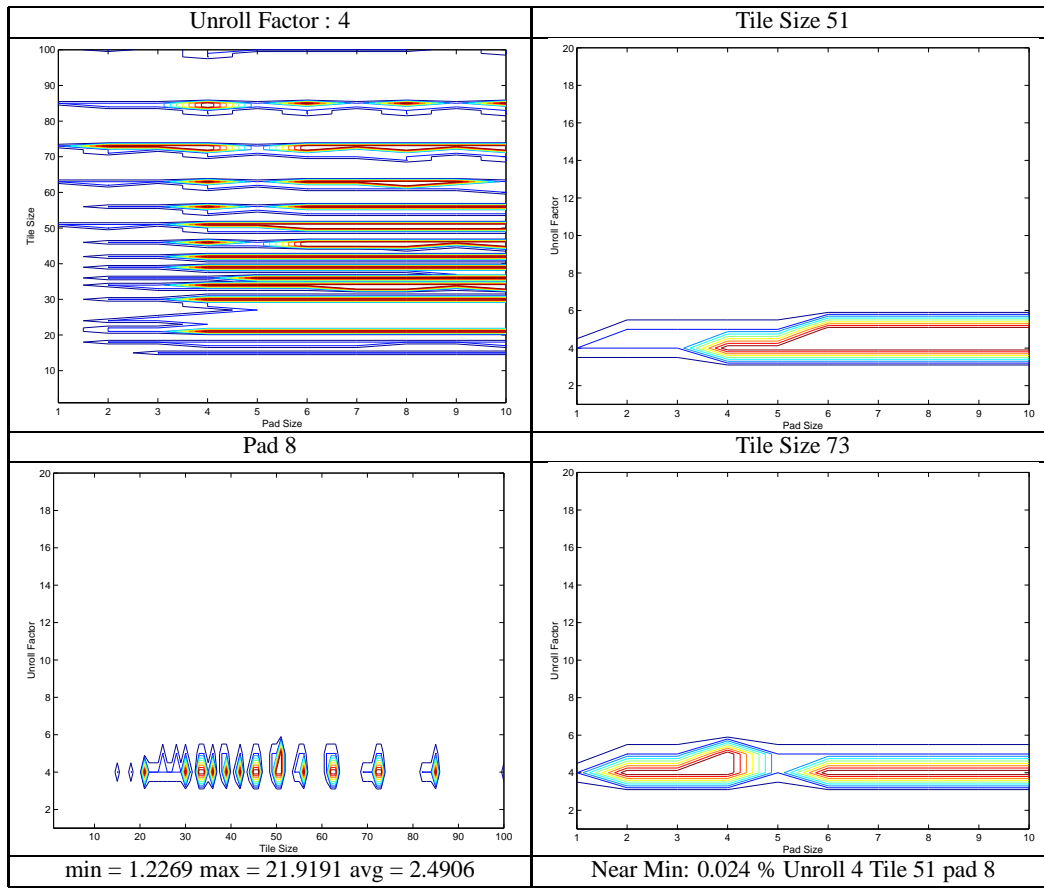


Figure 8. Transformation Space: UltraSparc N=512

has shown that a simple search algorithm can achieve good results. Future work will consider larger application programs and transformation spaces. The combination of static and iterative information to guide optimisation will also be considered. Finally, improved search algorithms will also be investigated now that we shown the usefulness of iterative compilation as a viable program optimisation approach, particularly for embedded systems.

Acknowledgment

We would like to thank Kyle Gallivan for his for his ideas and initial input to this work.

References

- [1] Aarts B., Barreteau, Bodin F., Brinkhaus P., Chamski Z., Charles H-P., Eisenbeis C., Gurd J.R., Hoogerbrugge J., Hu P., Jalby W., Knijnenberg P.M.W., O'Boyle M.F.P., Rohou E., Sakellariou R., Schepers H., Sez nec A., Stohr E.A. Verhoeven M., Wijshoff H.A.G. **OCEANS**, EuroPar97, Springer-Verlag, LNCS, Passau, August 1997.
- [2] S. Coleman and K. McKinley, **Tile Size Selection using Cache Organization and Data Layout**, Programming Language Design and Implementation, ACM Press, June 1995.
- [3] Hansen E., **Global Optimization Using Interval Analysis**, Marcel Dekker Inc. New York, 1992
- [4] Hummel S.F., Banicesu I., Wang C.-T. and Wein J., **Load Balancing and Data Locality via Fractiling: An Experimental Study**, Kluwer Academic Press, LCR, New York May 1995.
- [5] A. Nisbet., **GAPS: Genetic Algorithm Optimised Parallelisation** CPC 98, 7th Workshop on Compilers for Parallel Computing, Linkoepping, Sweden, June 1998.
- [6] Diniz P. and Rinard M., **Dynamic Feedback: An Effective Technique for Adaptive Computing**, Programming Languages Design and Implementation, Las Vegas 1997.
- [7] Fernandez M. , **Simple and effecient link-time optimizations of modula-2 programs** Programming Languages Design and Implementation, La Jolla, June 1995.
- [8] Bodin F., Chamski Z., Eisenbeis C., Rohou E. and Sez nec A. **GCDS: A Compiler Strategy for Trading Code Size Against Performance in Embedded Applications** IRISA Tech Rep. 1997.

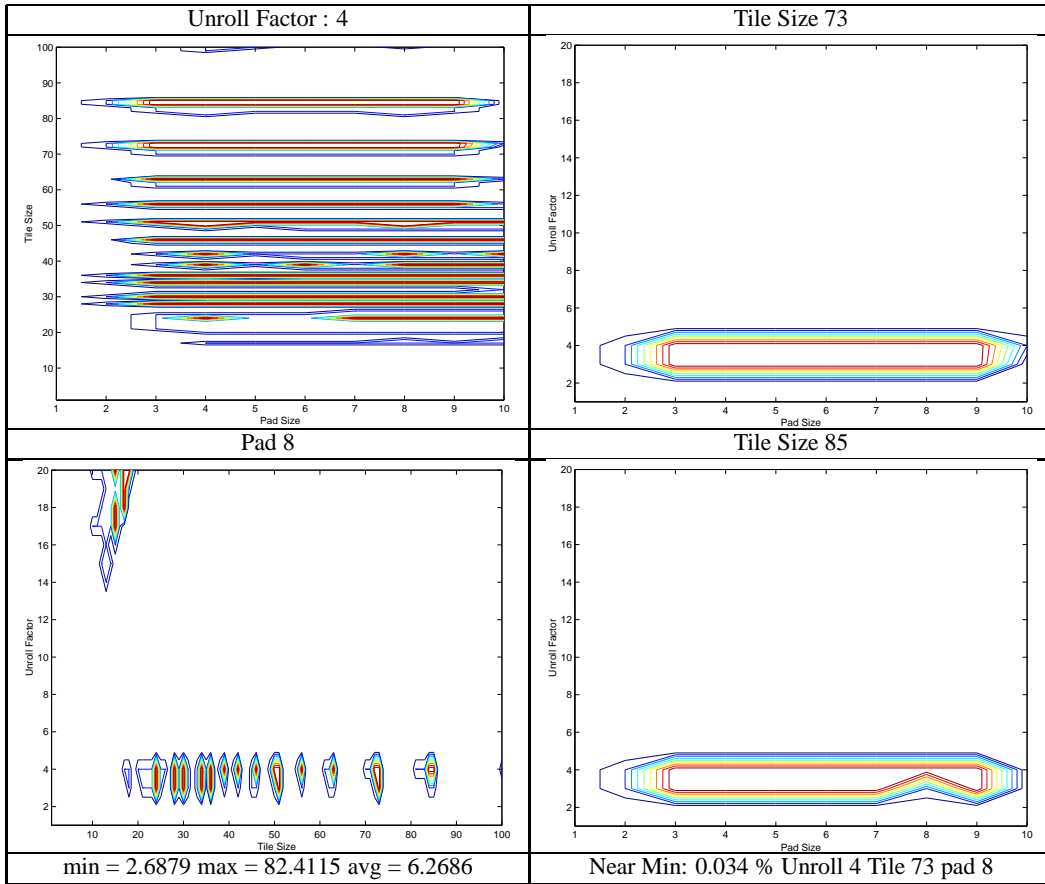


Figure 9. Transformation Space: Alpha N=512

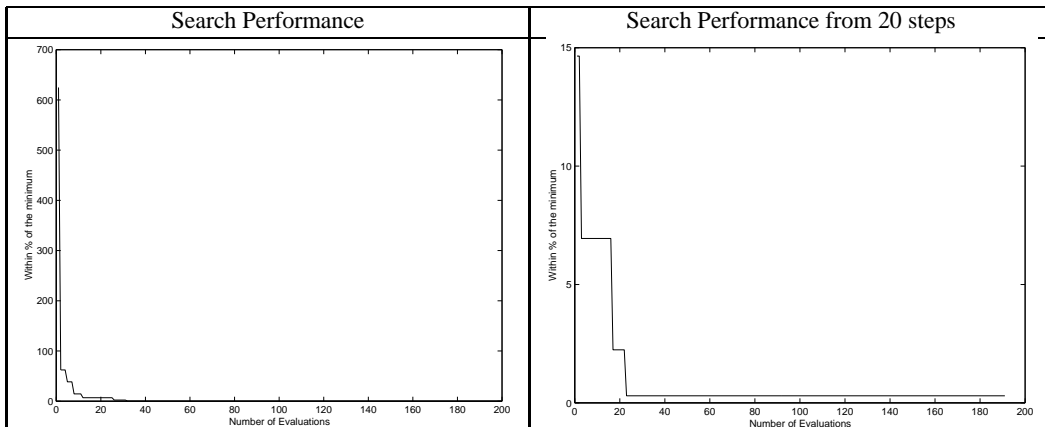


Figure 10. Performance of search algorithm : UltraSparc

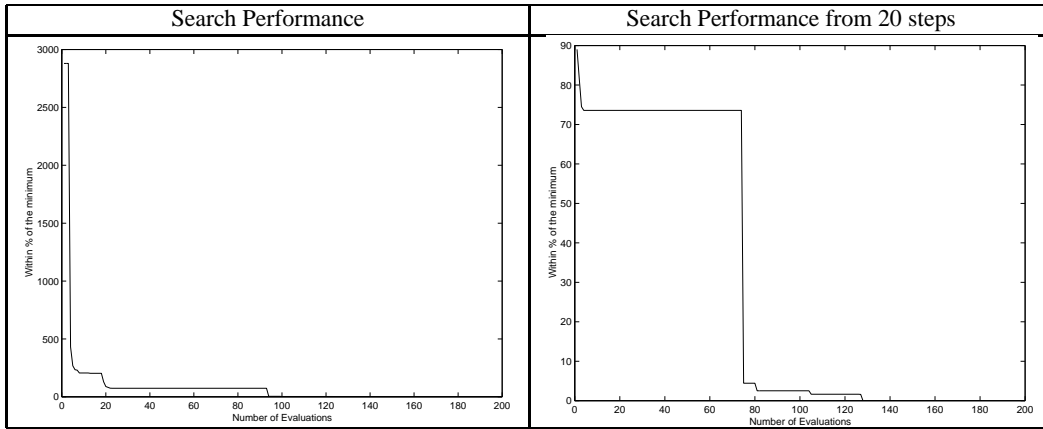


Figure 11. Performance of search algorithm : Alpha

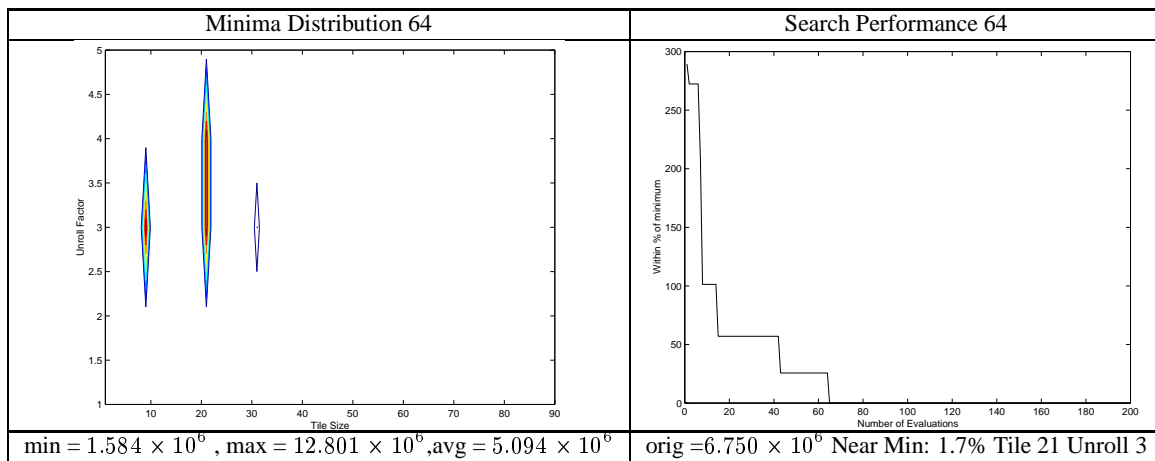


Figure 12. Minimal points and search algorithm performance for TM-1 N=64

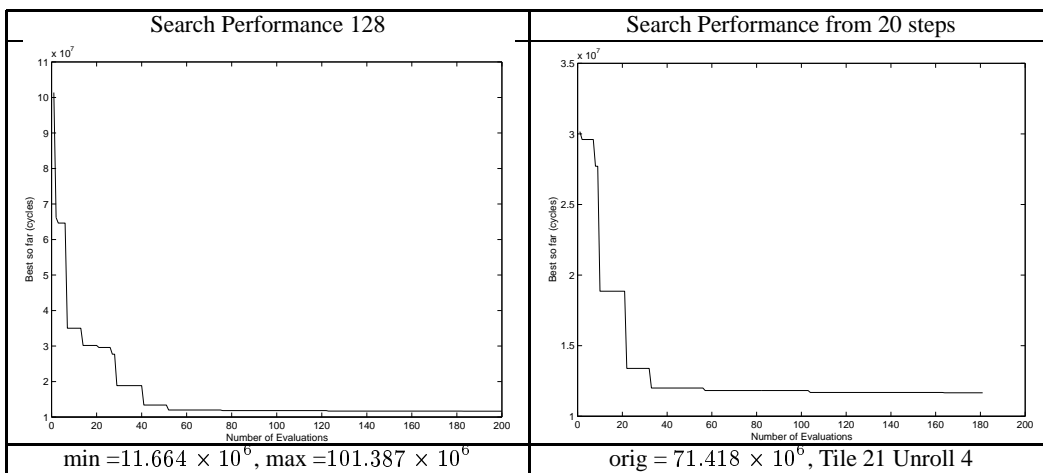


Figure 13. Minimal points and search algorithm performance for TM-1 N=128