

An Efficient Framework for Dynamic Reconfiguration of Instruction-Set Customization

Huynh Phung Huynh, Joon Edward Sim, Tulika Mitra
Department of Computer Science
National University of Singapore
Republic of Singapore 117543
{huynhph1, esim, tulika}@comp.nus.edu.sg

ABSTRACT

We present an efficient framework for dynamic reconfiguration of application-specific instruction-set customization. A key component of this framework is an iterative algorithm for temporal and spatial partitioning of the loop kernels. Our algorithm maximizes the performance gain of an application while taking into consideration the dynamic reconfiguration cost. It selects the appropriate custom instruction-sets for the loops and maps them into appropriate configurations. We model the temporal partitioning problem as a k -way graph partitioning problem. A dynamic programming based solution is used for the spatial partitioning. Comprehensive experimental results indicate that our iterative algorithm is highly scalable while producing optimal or near-optimal (99% of the optimal) performance gain.

Categories and Subject Descriptors

C.3 [Special-purpose and application-based systems]: Real-time and embedded systems; C.1.3 [Other Architecture Styles]: Adaptable architectures

General Terms

Algorithm, Performance, Design

Keywords

Customizable processors, instruction-set extensions, dynamic reconfiguration, temporal partitioning.

1. INTRODUCTION

Current generation embedded system designs are characterized by the increasing demand on higher performance under stringent time-to-market constraints. In this context, application-specific customizable processor cores strike the right balance between performance and design efforts. A customizable processor is, in general, configurable w.r.t. the micro-architectural parameters. More importantly, a customizable processor may support application-specific

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CASES'07, September 30–October 3, 2007, Salzburg, Austria.
Copyright 2007 ACM 978-1-59593-826-8/07/0009 ...\$5.00.

extensions of the core instruction-set. Custom instructions encapsulate the frequently occurring computation patterns in an application. They are implemented as custom functional units (CFU) in the datapath of the existing processor core. CFUs improve performance through parallelization and chaining of operations. Some examples of commercial customizable processors include Lx [13], ARCTM core [2], Xtensa [14] and Stretch S5 [3].

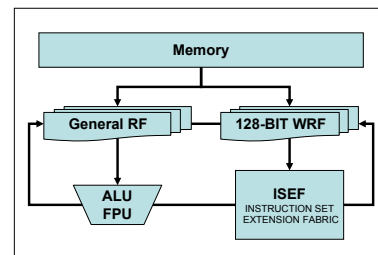


Figure 1: Stretch S5530 datapath.

Due to the limited area available for the CFUs in the datapath of the processor core, we may not be able to exploit all custom instruction enhancements of an application. Moreover, it may not be possible to increase the area allocated to the CFUs due to the linear increase in the cost of the associated system. In this context, runtime reconfiguration of the custom instruction-sets appears quite promising. Adapting to this trend, commercial customizable processors supporting dynamic reconfiguration have been proposed.

For example, Figure 1 shows the Stretch S5 engine [27] that incorporates Tensilica Xtensa RISC processor [14] and the Stretch Instruction Set Extension Fabric (ISEF). The ISEF is software-configurable datapath based on programmable logic. It consists of a plane of arithmetic/logic elements and a plane of multiplier elements embedded and interlinked in a programmable, hierarchical routing fabric. This configurable fabric acts as a functional unit to the processor. It is built into the processor's datapath, and resides alongside other traditional functional units such as the ALU and the floating point unit. The programmer defined application specific instructions (Extension Instructions) are implemented in this fabric. The core processor issues the Extension Instructions to ISEF, which performs the computation and returns the result. However, the distinguishing aspect of ISEF is that it is run-time configurable and reloadable. If the computation resource requirement of the custom instructions exceeds the capacity of ISEF, then the instructions can be partitioned into different configurations. When a user-defined instruction is issued, the S5 hardware checks to make sure the

corresponding configuration is loaded into the ISEF. If the required configuration is not present in the ISEF, it is automatically loaded prior to the execution of the user-defined instruction. In summary, the ISEF allows the system designers to define new instructions at runtime and thus extend the processor’s instruction-set.

However currently it is the programmer’s responsibility to manually choose and define the custom instructions and the configurations for architectures such as Stretch. Choosing an appropriate set of custom instructions for an application itself is a difficult problem. Significant research effort has been invested in developing automated selection techniques for custom instructions. Runtime reconfiguration has the additional complication of both *temporal and spatial partitioning* of the set of custom instructions in the reconfigurable fabric. In this paper, we develop a framework that starts with an application specified in ANSI-C and automatically selects appropriate custom instructions as well as clubs them into one or more configurations. A key component of our framework is an iterative algorithm for temporal partitioning of custom instructions based on k-way graph partitioning problem. A dynamic programming based pseudo-polynomial algorithm is used for the spatial partitioning of the custom instructions within a configuration. To the best of our knowledge, this is the first work that attempts automated custom instructions selection in the context of instruction-set extensible processor platforms with dynamic reconfiguration.

Even though most hardware-software partitioning solutions for FPGAs work at a coarse-grained level (such as task level and function level) and explore task-level parallelism [8, 19], we focus on *hot loop kernels* instead. This allows us to exploit instruction-level parallelism to significantly accelerate compute-intensive loops with custom instructions. Thus our framework first extracts a set of compute-intensive candidate loop kernels from the application. For each loop, we generate one or more custom instruction-set versions differing in performance gain and area tradeoffs in addition to the purely software version. The partitioning algorithm selects appropriate custom instruction-set versions for the loops implemented in fabric and clubs them into suitable configurations to achieve the highest performance gain.

Note that the reconfiguration cost model at task level [8, 19] and data flow graph level [26] are simple because the underlying directed acyclic graph representation ensures at most one reconfiguration between any two nodes. In contrast, our dynamic reconfiguration cost model is complex as the number of reconfigurations for one loop depends on temporal partitioning of all the other loops. Furthermore, our methodology allows custom instruction-sets corresponding to more than one loop to be placed within a single configuration. Thus spatial partitioning also plays a role in determining the performance gain of the application. The only other loop-level temporal partitioning work that we are aware of [23] considers one loop per configuration.

The remainder of this paper is structured as follows. Section 2 describes the system design flow. In Section 3, we present the problem formulation and a motivating example. Section 4 details our partitioning algorithm. Experimental setup and evaluation are described in Section 5. The related works are discussed in Section 6. Finally, Section 7 concludes the paper.

2. SYSTEM DESIGN FLOW

Figure 2 shows the system design flow. The input to the design flow is the C source code of the application we want to accelerate. The output is the custom-instructions accelerated application with synthesized datapath for each configuration.

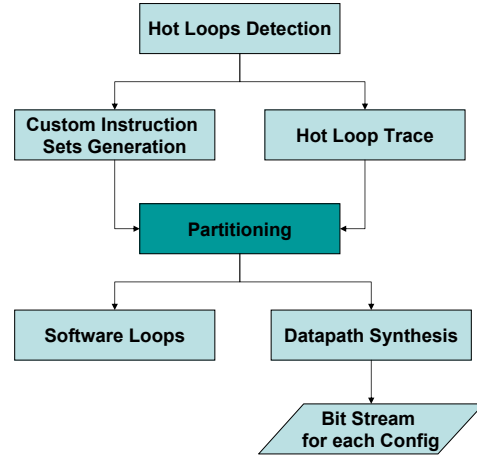


Figure 2: System design flow

Hot loops detection. Taking our cue from Amdahl’s law, we focus on the loops that take up a significant portion of the application’s total execution time. In particular, we define a loop with execution time greater than a certain percentage (typically $\geq 1\%$) of the application’s overall execution time to be a *hot* loop. The hot loop detector identifies such loops through profiling. Although the total number of loops in an application may be large, we consider only the hot loops to reduce the computation cost of the partitioning algorithm significantly. At the same time, the performance gain we obtain is still comparable to the case where all the loops of the application are considered. The cold loops may not increase the performance gain due to the additional reconfiguration overhead.

Custom instruction-set versions generation. Multiple custom instruction-set (CIS) versions are generated for each hot loop with a trade-off between hardware area and performance gain. A CIS version consists of a set of custom instructions extracted from the corresponding loop under an area constraint. Each CIS version is characterized by its area and performance gain. In general, the performance gain of a CIS version increases with larger area. To generate the CIS versions for a loop, the first step [4, 5, 10, 11, 18] identifies a large set of candidate patterns from the loop. Given this library of patterns, the second step selects a subset to maximize performance gain under hardware area constraint [4, 10, 9, 11, 22]. As the area increases, a CIS version with higher performance gain will be generated by selecting a larger subset. Moreover, different CIS versions can be generated by loop transformations such as loop unrolling, software pipelining, loop fusion, and others.

Loop Trace. The control flow among the hot loops is captured in the form of a loop trace (execution sequence of the loops) obtained through profiling. For typical embedded applications we have profiled, the number of hot loops and the loop trace size are quite small. For longer loop trace, we can use lossless compression techniques (such as SEQUITUR algorithm [25]) to compactly maintain the loop trace.

The hot loops with CIS versions and the loop trace are fed to the partitioning algorithm that decides the appropriate CIS version and configuration for each loop. The selected CIS versions to be implemented in hardware are then input into the datapath synthesis

tool. It generates the bit stream corresponding to each configuration (based on the result of temporal partitioning). These bitstreams are used to configure the fabric at runtime. The remaining loops are implemented in software on the core processor. Finally, the source code is modified to exploit the new custom instructions.

3. PROBLEM DEFINITION

We now formally define the partitioning problem for dynamic reconfiguration of instruction-set customization, which is the focus of this paper.

The input to the partitioning step is the set of hot loops $\{l_i | i = 1 \dots N\}$. Each loop is associated with multiple custom instruction-set (CIS) versions with a trade-off between hardware area and performance gain. Let $l_{i,j}$ (for $j = 1 \dots n_i$) be the j^{th} CIS version corresponding to loop l_i where n_i is the number of CIS versions of loop l_i . In addition, let $gain_{i,j}$ and $area_{i,j}$ denote the performance gain and area requirement of $l_{i,j}$. We assume that $l_{i,1}$ corresponds to the software loop without any custom instructions, i.e., $area_{i,1} = 0$ and $gain_{i,1} = 0$. For each loop l_i , only one of its CIS versions will be selected for implementation. For example, if $l_{i,1}$ is selected, loop l_i will be implemented in software without any custom instruction enhancements.

The control flow among the loop kernels is input in the form of a loop trace. Finally, $MaxA$ represents the hardware area available for each configuration and ρ represents the time required for a single reconfiguration. We assume partial reconfiguration is not supported, i.e., a configuration is completely replaced by another configuration in the fabric. Hence both $MaxA$ and ρ are constants.

We do not allow intra-loop reconfiguration to avoid high reconfiguration cost. Thus the custom instructions corresponding to a loop cannot straddle across configuration boundaries. In other words, the selected CIS version of a loop is completely accommodated within a configuration. Each configuration, however, consists of CIS versions corresponding to one or more loops. Thus the problem boils down to

1. *Temporal partitioning* of the loops selected for hardware acceleration with CIS into one or more configurations, and
2. *Spatial partitioning* of the loops within a configuration by selecting appropriate CIS version for each loop

The performance gain of the application (PGA) is then defined as

$$PGA = \left(\sum_{i=1}^N \sum_{j=1}^{n_i} s_{i,j} \times gain_{i,j} \right) - r * \rho \quad (1)$$

$$\sum_{j=1}^{n_i} s_{i,j} \leq 1 \quad (2)$$

where r is the number of reconfigurations given the partitioning and $s_{i,j}$ is a binary variable equal to 1 if CIS version $l_{i,j}$ is selected and 0 otherwise.

Dynamic reconfiguration through temporal partitioning enlarges the available area for the design by increasing the number of configurations. Therefore, each loop can select better CIS version to be implemented in hardware and better performance gain will be achieved. However, this increase in number of configurations may not result in better overall performance due to the reconfiguration cost. On the other hand, if we minimize the number of configurations, the available area is quite restricted. Consequently, each loop will select its CIS version with smaller area and the performance gain of the application is much smaller, especially when the reconfiguration cost is smaller.

Our objective is to maximize the performance gain by selecting an appropriate CIS version for each loop and mapping it into an appropriate configuration.

3.1 Motivating Example

Loop	Version	Area (#AU)	Gain (K cycles)
loop1	1	0	0
	2	257	111
	3	301	160
	4	1612	563
loop2	1	0	0
	2	761	230
	3	1041	387
	4	1321	426
	5	2004	556
loop3	1	0	0
	2	967	493
	3	1249	549

Table 1: CIS versions for 3 loops in our motivating example.

For our motivating example, we consider an application with three hot loops: `loop1`, `loop2` and `loop3`. Table 1 shows the performance/silicon area tradeoff of different custom instruction-set versions for each loop. In particular, the table shows the hardware requirement in terms of arithmetic units (AU) and corresponding performance gain in terms of K cycles. For example, `loop3` has three CIS versions. Version 1 of each loop is the software version (without any custom instructions enhancements) with zero area and performance gain. We need to select appropriate CIS versions for the three loops under hardware area constraint for a configuration of 2048 AUs. The cost for a single reconfiguration is 15K cycles. The graph on the left-hand side of Figure 3 shows control flow information among the loops for this example. The actual input to our algorithm is the loop trace. We use the graph here (derived from the loop trace) for illustration purposes. We will, however, use a similar graph (called reconfiguration cost graph) later in our temporal partitioning algorithm.

If the system does not support dynamic reconfiguration, the best solution (solution (A) in Figure 3) under the hardware area constraint is the selection of version 3 of `loop1`, version 2 of `loop2`, and version 2 of `loop3`. Total performance gain is $160 + 230 + 493 = 883$ K cycles and there is no reconfiguration cost.

However, in the presence of dynamic reconfiguration, we can improve the solution. A trivial solution is to put each loop into one configuration (solution (B) in Figure 3). We can then select the CIS version of a loop with the largest area less than or equal to the area of a configuration: version 4 for `loop1`, version 5 for `loop2` and version 3 for `loop3`. Total performance gain is $563 + 556 + 549 = 1668$ K cycles and the total reconfiguration cost is $(20 + 11 + 9 + 9) \times 15 = 735$ K cycles. Therefore the resulting net performance gain after subtracting the reconfiguration cost is $1668 - 735 = 933$ K cycles. While the net performance gain is better than the case when dynamic configuration is not supported, it is not the optimal solution.

The optimal solution is to put `loop2` and `loop3` into one configuration and `loop1` into a different configuration (solution (C) in Figure 3). CIS versions 4, 3, and 2 will be selected for `loop1`, `loop2`, and `loop3`, respectively. The performance gain is 1443 K cycles while reconfiguration cost is $(9 + 9) \times 15 = 270$ K cycles. Hence, the net performance gain is $1443 - 270 = 1173$ K cycles.

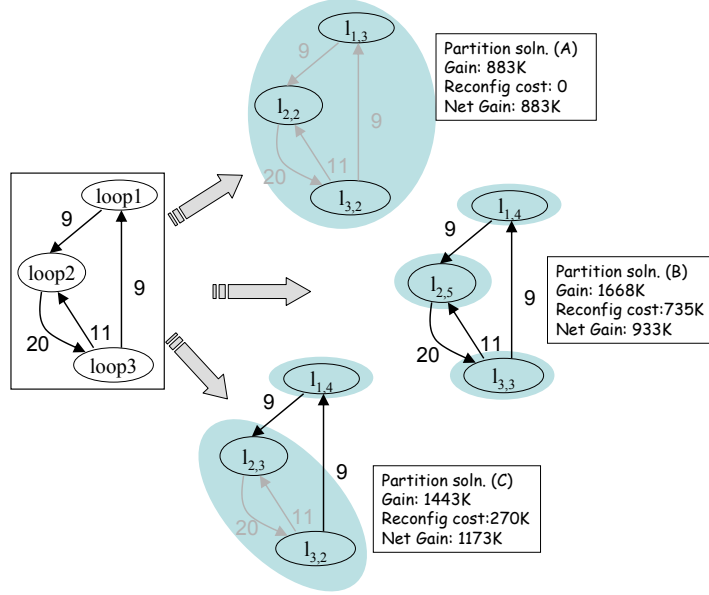


Figure 3: Some partitioning solutions for our motivating example.

4. PARTITIONING ALGORITHM

Finding the optimal combination of temporal and spatial partitioning is a difficult problem. Given N loops, the number of possible configurations is 2^N . However, the number of ways to partition N loops into mutually-exclusive configurations corresponds to the $N + 1^{th}$ Bell number. According to de Bruijn [12], asymptotic limits of Bell numbers is $O(e^{N \ln(N)})$.

Our partitioning algorithm needs to make three choices: (1) optimal number of configurations k , (2) temporal partitioning of the loop kernels into k configurations, and (3) spatial partitioning of the loop kernels in each configuration, i.e., choosing the appropriate custom-instruction set (CIS) version for each loop kernel. Clearly, these choices are inter-dependent. The selection of CIS versions for the loops determines the partitioning solution and vice versa.

4.1 Overview

Algorithm 1: Iterative Partitioning Algorithm

Input: Set of hot loops with custom instruction-set versions: L
Loop Trace: T
Maximum Area of a configuration: $MaxA$
Reconfiguration Cost: ρ

Result: Partition with the best net performance gain

for $k = 1$ to $|L|$ **in steps of 1** **do**

```

   $C := \text{global\_spatial\_partition}(L, k \times MaxA);$ 
   $P := \text{temporal\_partition\_with\_CIS}(C, T, k);$ 
   $P' := \text{temporal\_partition\_wo\_CIS}(T, k);$ 
   $soln := \text{local\_spatial\_partition}(L, P, MaxA);$ 
   $soln' := \text{local\_spatial\_partition}(L, P', MaxA);$ 
  if  $\text{net\_gain}(soln') > \text{net\_gain}(soln)$  then  $soln := soln'$ ;
  if  $\text{net\_gain}(soln) > \text{net\_gain}(bestSoln)$  then  $bestSoln := soln$ ;

```

end

return $bestSoln$;

We propose an iterative algorithm (Algorithm 1) for joint temporal and spatial partitioning of the custom instruction-sets corresponding to the hot loop kernels. The algorithm iterates from a con-

straint of having exactly 1 configuration (i.e., no reconfiguration) to the upper bound of having $|L|$ configurations where L is the set of hot loops. The solutions (A) and (B) in our motivating example (see Figure 3) represent the two extremes ($k = 1$ and $k = |L|$) while the remaining iterations explore the rest of the design space.

For the iteration with k configurations, we would like to identify the k -way partitioning solution with the optimal net performance gain. Unfortunately, temporal and spatial partitioning are again dependent on each other due to the reconfiguration cost. To break this cycle, we apply a heuristic technique. The heuristic first assumes that we have a continuous area of $k \times MaxA$ available to us where $MaxA$ is the maximum area for a configuration. The assumption of continuous area allows us to tentatively select optimal CIS versions for the loops in an ideal (but un-realizable) situation where reconfiguration cost is zero. In reality, we have k distinct configurations with $MaxA$ area each. So we partition the loop kernels with selected CIS versions into k configurations such that each configuration has *roughly* $MaxA$ area and the reconfiguration cost is minimized. As we break up the continuous area into k distinct areas, some configurations end up being bigger than $MaxA$, while some other configurations are smaller than $MaxA$. To fix this, we have a final patch-up stage that performs spatial partitioning within each configuration to re-distribute $MaxA$ space among the constituent loop kernels.

Figure 4 illustrates the three phases of the iterative algorithm corresponding to the iteration with number of configurations equals to 2. The input is the three loops in the motivating example and their CIS versions. For example, `loop1` has 4 CIS versions $l_{1,1}$, $l_{1,2}$, $l_{1,3}$ and $l_{1,4}$ in order of increasing area and performance gain.

The first phase **global_spatial_partition** partitions the area $k \times MaxA$ (where k is the number of configurations for that iteration) among the loops by selecting the CIS versions such that the performance gain is optimal. This phase disregards the reconfiguration cost. It also assumes that a continuous hardware area of size $k \times MaxA$ is available for hardware acceleration of all the loops. We have developed a *Dynamic Programming* algorithm for this phase.

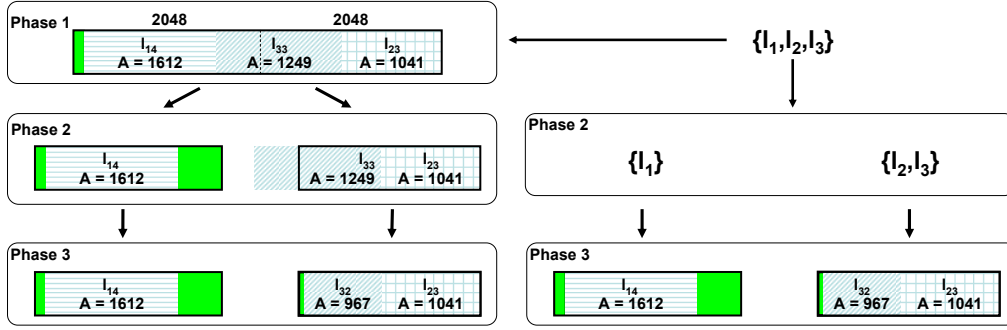


Figure 4: Three phases of iterative partitioning algorithm for number of configurations = 2

This phase may choose to select the software version for some loops. For our running example, the first phase in Figure 4 chooses CIS versions $l_{1,4}, l_{2,3}, l_{3,3}$ when the area budget equals to the area for 2 configurations.

After the first phase, we have the set of selected CIS versions C for the hot loops. However, we cannot implement this solution as (1) the reconfiguration cost has not been considered and (2) the loops still need to be partitioned into different configurations. In the second phase **temporal partition with CIS**, we perform temporal partitioning of the hot loops into k configurations such that the reconfiguration cost is minimized. This phase returns the partitioning solution P for the set of loops selected for custom instructions enhancements from the first phase. In this phase we also find an alternative partitioning solution P' for the original set of hot loops, i.e., it disregards the results of the first phase. This partitioning **temporal partition wo CIS** only considers the reconfiguration cost and ignores the CIS versions. Partition P gives good results when performance gain of CIS versions is high relative to the reconfiguration cost. On the other hand, partition P' gives better results for the case when the reconfiguration cost is high relative to the performance gain. P and P' complement each other in the search for the best partitioning solution. We model the temporal partitioning as k -way *weighted graph partitioning problem*, which is well studied [16, 17].

In Figure 4, the left hand side shows the partition P and the right hand side shows the partition P' . For P , the second phase partitions the three loops with selected CIS versions into two configurations: $l_{1,4}$ in the first configuration and $l_{2,3}, l_{3,3}$ in the second configuration. On the other hand, P' simply partitions the three loops based on reconfiguration cost into two configurations. In this example, P and P' return the same temporal partitioning. However, due to the reconfiguration cost, P and P' may be different.

We now have k configurations for each partitioning solution P and P' . The k -way weighted graph partitioning produces partitions with roughly equal size. Therefore for partition P , the area requirement of some of the configurations may exceed the maximum area $MaxA$. Partitioning solution P' , on the other hand, does not select any CIS version a-priori. Thus, for each configuration in P and P' , the third phase **local spatial partition** locally selects the CIS versions for the loops in that configuration to maximize performance gain under area constraint $MaxA$. We again use dynamic programming to perform optimal spatial partitioning for each configuration.

In Figure 4, for partition P , the area requirement of the second configuration exceeds the maximum area budget. Hence phase

3 for this partition replaces CIS version $l_{3,3}$ with $l_{3,2}$. Phase 3 keeps the CIS version for loop l_1 unchanged even though there is additional area available (the green part) as $l_{1,4}$ is the best version for l_1 . However, in general, the additional area can lead to the selection of better versions for some loops. Typically, the sum of the performance gain of all the loops in phase 3 is at least 90% of the performance gain of phase 1. The third phase of P' simply selects CIS versions of the loops in each configuration for the first time.

If the net performance gain of the current solution is better than the best solution so far, we update the best solution. Then we start a new iteration with $k = k + 1$. The algorithm terminates when in the current solution, each loop has been assigned its CIS version with the best performance gain. In the worst case, the algorithm runs for $|L|$ iterations. With the motivating example, our algorithm returns the optimal solution, which has two configurations (see Figure 4) and the performance gain is 1173K cycles.

4.2 Spatial Partitioning

We propose a pseudo-polynomial time dynamic programming algorithm to select the appropriate CIS versions for the loops such that the performance gain is optimal under a hardware area budget. This algorithm is employed in the first phase and the third phase of our iterative solution with different parameters.

Let $G_i(A)$ be the *maximum* performance gain of loops $l_1 \dots l_i$ under an area budget A . Then $G_i(A)$ can be defined recursively.

$$G_i(A) = \max_{\substack{j=1 \dots n_i \\ area_{i,j} \leq A}} (gain_{i,j} + G_{i-1}(A - area_{i,j})) \quad (3)$$

That is, given an area A , we explore all possible CIS versions for l_i and choose the one that results in maximum performance gain for loops $l_1 \dots l_i$. The base case for loop l_1 is

$$G_1(A) = \max_{\substack{j=1 \dots n_1 \\ area_{1,j} \leq A}} (gain_{1,j}) \quad (4)$$

The maximum performance gain for loops $l_1 \dots l_N$ under area budget $AREA$ then corresponds to $G_N(AREA)$.

Algorithm 2 encodes this recursion as a bottom-up dynamical programming algorithm. The step value Δ determines the granularity of area. It is chosen based on the minimum area difference between two successive CIS versions for any loop. The time complexity of this algorithm is $O(N \times \frac{Area}{\Delta} \times x)$ where $x = \max_{i=1 \dots N} (n_i)$.

Algorithm 2: Spatial Partitioning

Input: Set of interesting loops l_1, l_1, \dots, l_N with CIS versions;
Area constraint: $AREA$
Result: Maximum performance gain
for $A = 0$ to $AREA$ in steps of Δ **do**
 $G_1(A) \leftarrow \max_{\substack{j=1 \dots n_1 \\ area_{1,j} \leq A}} (gain_{1,j})$
end
for $A = 0$ to $AREA$ in steps of Δ **do**
 for $i=2$ to N **do**
 $G_i(A) \leftarrow \max_{\substack{j=1 \dots n_i \\ area_{i,j} \leq A}} (gain_{i,j} + G_{i-1}(\lfloor \frac{A - area_{i,j}}{\Delta} \rfloor \times \Delta))$
 end
return $G_N(AREA)$;

4.3 Temporal Partitioning

We map our temporal partitioning problem to k-way weighted graph partitioning problem. The k-way weighted graph partitioning problem is defined as follows: Given an undirected graph $G = (V, E)$ with weights both on the vertices and the edges, partition V into k subsets V_1, V_2, \dots, V_k such that $V_i \cap V_j = \emptyset$ for $i \neq j$, $\bigcup_i V_i = V$, the sum of the vertex-weights in each subset is roughly equal and the sum of the edge-weights whose incident vertices belong to different subsets (edge-cut weights) is minimized.

We generate a **Reconfiguration Cost Graph (RCG)** from the loop trace for k-way weighted graph partitioning. After the first phase, we have tentatively selected CIS versions for the loops. Each vertex in the RCG represents a hot loop selected for hardware acceleration in the first phase. In other words, we do not consider the loops for which the first phase selects software-only version. Given a vertex v associated with loop l , we assign the area of the CIS version selected for l as the weight of the vertex v . When CIS versions from the first phase are ignored, the RCG includes all the loops and we assume unit hardware cost for each vertex.

The edge weight between vertex v (corresponding to loop l) and v' (corresponding to loop l') is defined as the reconfiguration cost between loop l and loop l' if they are mapped to two different configurations. The edge between v and v' exists if and only if control can flow from loop l to l' or l' to l without passing through any other hot loops. The weight on the edge between v and v' represents the number of times control flows directly from loop l to l' and l' to l multiplied by the single reconfiguration cost ρ . This weight can be derived from the loop trace as follows. If we eliminate the software-only loops from the loop trace, then the weight is the the number of times the string ll' and $l'l$ appear in the loop trace multiplied by ρ . The time complexity of creating RCG is linear in the size of the hot loop trace.

Figure 5 shows an example of RCG generation from the loop trace. It shows a loop trace $ABCBCBA$ of three hot loops A, B, C . The reconfiguration cost ρ is assumed to be 1 time unit. If all the loops are selected to be placed in hardware, then there are 2 reconfiguration points between loops A and B if they are partitioned into different configurations. Similarly, there are 4 reconfiguration points between loops B and C if they are partitioned into different configurations. However, there are no reconfiguration points between loops A and C directly as the control transfers between them always pass through B . However, if we choose to implement B in software in the first phase, then B is eliminated from the RCG. In this case, there are 2 reconfiguration points between loops A and C if they are partitioned into different configurations.

The objective now is to partition the RCG into k configurations such that the configurations have roughly equal area (or the configurations have roughly equal number of loops when area is ignored)

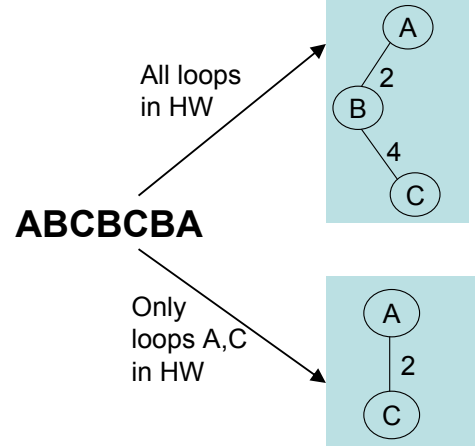


Figure 5: Reconfiguration cost graph from loop trace

and the reconfiguration cost (edge-cut weights) is minimized. If the configurations have roughly equal area, then the loops have higher probability of retaining the optimal CIS versions selected in the first phase regardless of the third phase. As a result, total performance gain (excluding reconfiguration cost) after the third phase is near the optimal performance gain in the first phase. The rationale behind having roughly equal number of loops in each configuration when CIS versions are ignored (by assigning unit cost to each vertex in the RCG), is to create a balanced temporal partition. It ensures that equal number of loops compete for each configuration space during subsequent spatial partitioning.

We use multilevel k-way partitioning scheme by Karypis and Kumar [17]. The multilevel partitioning scheme consists of three phases: coarsening phase, partitioning phase and uncoarsening phase. During coarsening phase, a sequence of smaller graphs $G_i = (V_i, E_i)$, each with fewer vertices, is constructed from the original graph $G_0 = (V_0, E_0)$ such that $|V_i| < |V_{i-1}|$. The coarsening phase ends when the coarsest graph G_m has a small number of vertices or the reduction in the size of successively coarser graph becomes too small. Then, the partitioning phase computes a k-way partitioning P_m of the coarse graph $G_m = (V_m, E_m)$ such that each partition contains roughly $|V_0|/k$ vertex weight of the original graph. The k-way partitioning of G_m is computed using multilevel bisection algorithm [16]. During the uncoarsening phase, the partitioning P_m of the coarser graph G_m is projected back to the original graph by going through the graphs $G_{(m-1)}, G_{(m-2)}, \dots, G_1$. At each intermediate level, the partitioning is refined based on Kernighan-Lin [20] partitioning algorithm and their variants.

5. EXPERIMENTAL EVALUATION

We have developed two algorithms, exhaustive search and greedy search, for the purpose of evaluating our proposed methodology. The results of the two algorithms are compared with our proposed methodology in two different sets of experiments. In the first set of experiments, we run the three algorithms using synthetic input to evaluate the scalability and efficiency of the algorithms. We generate input data with 5 to 100 hot loops for this set of experiments. In the second set of experiments, we conduct a case study of the JPEG application with custom instructions implemented on a commercial platform Stretch 5 [27] that supports runtime reconfiguration.

Exhaustive Search. The exhaustive search algorithm computes the optimal results by evaluating all possible temporal and spatial partitioning. We use the algorithm described in Kreher and Stinson [21] to enumerate all possible partitions. We then find the optimal implementation of each configuration in the partition by choosing CIS versions of the constituent loops through our spatial partitioning algorithm. The net gain of each enumerated partition is then estimated by a brute force computation of the reconfiguration cost by traversing the loop trace. The partition with the maximum net performance gain is then the optimal solution. Our experiments show that the exhaustive search algorithm cannot scale with increasing number of hot loops.

Algorithm 3: Greedy Search Algorithm

Input: Set of hot loops with custom instructions: L
 Loop Trace: T
 Maximum Area of a configuration: $MaxA$
 Reconfiguration Cost: ρ

Result: Partitioning solution

```

current := new_configuration();
continue := true;
while continue = true do
  C := compute_reconfig_cost_for_unselected_loops(L);
  li,j :=
  select_most_profitable_feasible_CIS(C, L, T, MaxA, solution);
  if li,j is not found then
    if current is not empty then
      update solution by adding current;
      current := new_configuration();
    else
      continue := false;
    end
  else
    update current with li,j;
    remove from L all CIS versions of loop li;
  end
end
return solution

```

Greedy Search. The greedy search algorithm (see Algorithm 3) constructs a solution by building one configuration at a time until no more CIS version can be added without causing a degradation in performance. The input is the set of hot loops with custom instruction-set versions L , loop trace T , area constraint $MaxA$, and single reconfiguration cost ρ . A solution consists of one or more configurations. The algorithm begins with an empty solution and an empty current configuration.

In each iteration, we pre-compute a reconfiguration cost array C . For any unselected loop l_i , the array C gives the expected additional reconfiguration cost if l_i is added to the current configuration. Given C , the current solution and the current configuration, we can now compute the expected performance gain of each CIS version if we add it to the current configuration. For CIS version $l_{i,j}$, this expected performance gain is estimated by subtracting from $gain_{i,j}$, the additional reconfiguration cost for loop l_i (available from array C). We now select the CIS version with the maximum expected *positive* performance gain that can be added to the current configuration without violating the area constraint. The selected CIS version is then added to the current configuration. All the other CIS versions of the same loop are subsequently removed from the set L .

In the event that no CIS version can be selected, there are two possibilities. The first possibility is that no more loops can be added to the current configuration without violating the area constraint (*current* configuration is not empty in Algorithm 3). In this case, we update the solution with the current configuration and re-start

the process of selecting CIS versions with an empty configuration. The second possibility is that no more loops can be added to the current solution without decreasing its net performance gain (*current* configuration is empty, i.e., we are trying to select the CIS version under maximum area constraint). In this case, the algorithm stops and returns the solution built so far.

5.1 Efficiency and Scalability of Algorithms

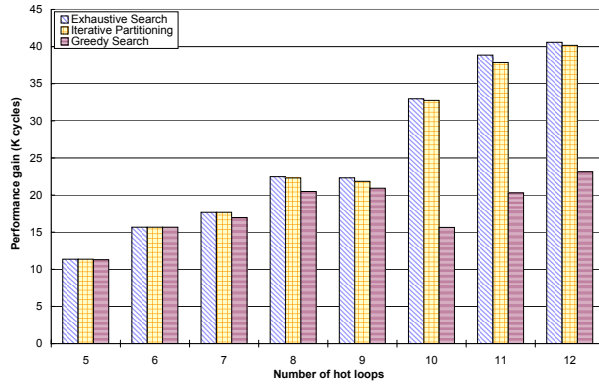
For this experiment, we generate synthetic inputs with number of hot loops ranging from 5 to 100. The number of CIS versions for each loop is generated randomly and ranges between 1 to 10. The performance gain of each CIS version ranges between 1000 to 10,000 time units. The hardware area is between 1 to 100 units. The performance gain increases with hardware area for each loop.

The reconfiguration costs between two loops, if they are assigned to different configurations, are generated randomly. They are in the range 0 to $maxCost$ where $maxCost$ is approximately 40-50% of the average performance gain of all the CIS versions of all the loops $\frac{\sum_{i=1}^N \sum_{j=1}^{n_i} gain_{i,j}}{\sum_{i=1}^N n_i}$. The value of $maxCost$ ensures that the reconfiguration cost is neither too high nor too low. Both the extremes reduce the search space considerably. If the reconfiguration cost is too high, we should only consider partitions with a small number of configurations. If the reconfiguration cost is too low, then the solution is to simply select the CIS version with the highest speedup for each loop and construct as many configurations as required. The hardware area constraint $MaxA$ is approximately 20-30% of the sum of the average area requirements of the CIS versions of all the loops $\sum_{i=1}^N \frac{\sum_{j=1}^{n_i} area_{i,j}}{n_i}$. This ensures that all the loops with their CIS versions cannot fit under the area constraint.

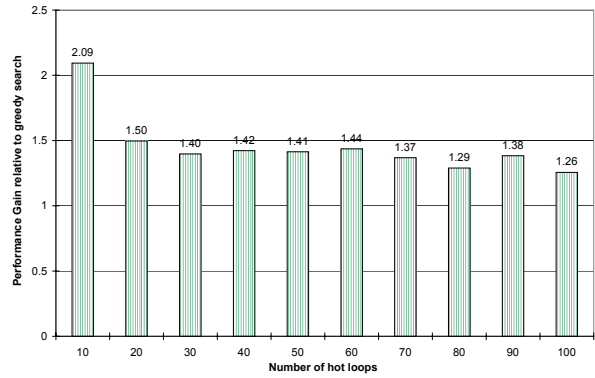
Number of Hot Loops	Running time (sec)		
	Exhaustive search	Greedy search	Iterative partitioning
5	0.26	0.01	0.07
6	1.34	0.02	0.07
7	7.84	0.01	0.07
8	43.91	0.01	0.09
9	283.22	0.04	0.07
10	1788.20	0.01	0.11
11	12604.33	0.01	0.13
12	86338.37	0.01	0.15
20	N.A.	0.02	0.48
40	N.A.	0.04	4.30
60	N.A.	0.07	18.25
80	N.A.	0.11	55.61
100	N.A.	0.16	118.76

Table 2: Running time of the algorithms for synthetic input.

Table 2 shows the running times of the three algorithms for synthetic input with different number of hot loops. The running time of the exhaustive search algorithm, while relatively small with smaller number of loops, increases by almost an order of magnitude each time one more loop is considered. The results of exhaustive search for more than 12 loops cannot be obtained even after waiting for a day. On the other hand, although iterative partitioning algorithm is slower than greedy search in general, its running time is quite acceptable (less than 2 minutes). This demonstrates the scalability of our approach. Moreover, iterative partitioning generates much better quality solutions compared to greedy search as presented in the following.



(a) Comparison of the performance gain of the algorithms for input with 5-12 hot loops



(b) Relative performance gain of iterative partitioning compared to greedy search

Figure 6: Experiment results for synthetic benchmarks

Figure 6(a) and 6(b) compare the quality of the solutions returned by the three different algorithms with number of hot loops varying from 5 to 12. Figure 6(a) shows that the performance gain obtained using our approach is close to the optimal gain obtained with exhaustive search while greedy search falls far behind.

Figure 6(b) presents the comparison between the performance gain of iterative partitioning and greedy search for input with more than 12 hot loops (where exhaustive search fails to produce results). The iterative algorithm consistently outperforms greedy search by a factor of 1.26 to 2.09.

5.2 Case Study of JPEG

Loop ID	#AU	#MU	Gain (K cycles)
0	2249	4096	32
1	1612	2880	563
	257	704	111
	389	2176	254
2	2004	6272	556
	1041	2048	387
	1321	2592	426
	761	1504	230
3	207	0	493
	424	2	549
4	2515	1536	1094
5	1530	3584	1669
	1300	3584	1643
6	981	4480	1095
	491	2240	739
	393	1792	590
7	1059	2880	511
8	1089	2880	91
9	1764	1280	194
	1114	768	188

Table 3: CIS versions for JPEG. The area requirements are in terms of arithmetic units (AU) and multiplier units (MU).

We present a case study of the JPEG image compression algorithm. In this study, we envision a scenario in which an image is decoded and then encoded subsequently. The hot loops are profiled and the loop trace is generated using an in-house tool based on

OpenImpact[1], an open source compiler. The profiling works in two phases. The timing information of each loop is collected by inserting appropriate time stamps at the entry and exit points of the loops. After the first pass, loops which take up more than 1% of the computation time can be detected. During the second pass, the compiler inserts appropriate code to capture the entry point of the hot loops. The resulting application, when executed, generates a trace of the hot loops.

Our loop profiler is able to identify more than 15 hot loops. For our experimental purposes, we have selected 10 loops for which custom instruction versions are manually generated for the Stretch S5 platform [27]. The profiler in Stretch IDE can then provide us the performance gain and hardware area of the CIS versions of each loop. Table 3 shows the various CIS versions for each loop and their respective area requirements and performance gain. It is worth noting that the performance gain of the CIS versions do not commensurate with area increase in general. For example, loop 0 takes up 2249 arithmetic units and 4096 multiplier units but only gives 32K cycles of performance gain. In contrast, the CIS versions of loop 3 use far less area but give much better performance. This is because the parallelism that can be exploited varies from one loop to another.

The configuration time of the whole fabric of Stretch development board, which includes 4096 4-bit arithmetic units (AUs) and 8192 4-bit \times 8-bit multiplier units (MUs) is approximately 100 μ s. Given that the CPU runs at 300MHz, the configuration time translates to roughly 30K CPU cycles. We define *one hardware area unit* to be a tuple of 400 AUs and 800 MUs. Since the configuration time is proportional to the size of the fabric, configuration time of one hardware area unit is approximately 3K CPU cycles. By scaling the configuration time according to the fabric size, we can easily compute the configuration time for any fabric size.

It is possible to fit CIS versions of all the hot loops from our JPEG application in a suitably-sized fabric. For our experimental purposes, we assume that the hardware area constraint varies from one hardware area unit to 20-30% of the sum of maximum hardware area for all the loops (5 – 15 hardware units for JPEG application). This will lead to the necessity of dynamic reconfiguration. We run all the three algorithms (Exhaustive search, Greedy search and Iterative partitioning) under these different area constraints. Our profiling data indicates that the application takes up around 20 million cycles on Stretch CPU without custom instructions enhancements. It should be noted, however, that the speedup we

obtain for a particular application depends on the quality of the custom instructions generated in the first place. Our focus in this experiment is to evaluate our proposed algorithm in comparison with Greedy search and Exhaustive search. That is, we are only concerned about comparing the performance gain obtained using the different algorithms starting with the same set of CIS versions. Our results show that the our proposed algorithm is always optimal or near-optimal and produces much better results than Greedy search most of the time.

In Figure 7(a), we evaluate the performance gain possible if dynamic reconfiguration is exploited. We compare the performance gain obtained using Iterative partitioning and Greedy search with the case when no reconfiguration is allowed. Clearly, Iterative partitioning and Greedy search can choose to use more than one configuration. However, the algorithm for no reconfiguration case is restricted to a single configuration and hence only performs spatial partitioning. For each hardware area unit, the left most column represents the performance gain obtained by our algorithm while the central and the rightmost column represent the performance gain under Greedy search and the single configuration, respectively.

If we compare the results of our algorithm with that of only one configuration, the advantage of exploiting dynamic reconfiguration decreases as the hardware area increases. This is to be expected, as more custom instructions can fit into the larger area to gain suitable speedup, thus reducing the need to virtualize hardware through runtime reconfiguration. On the other hand, the graph demonstrates that our algorithm increases the performance gain over and above single configuration by at least 34% and as much as 78%.

However, the true strength of our algorithm is not demonstrated by comparing results with no reconfiguration case. The Greedy search algorithm demonstrates that a simple heuristic fails to achieve substantial performance gain over no reconfiguration case. The Greedy search algorithm fails to capitalize dynamic reconfiguration as much as our algorithm. Often the Greedy search performs as good as the single configuration, and in some cases, even worse. On the other hand, our proposed methodology always performs better than the Greedy search, being at least 14% and as much as 91% better than Greedy search.

Figure 7(b) measures how closely our proposed methodology approximates the optimal results obtained through Exhaustive search. The graph shows that our algorithm returns solution that coincides with the optimal solution most of the time, while falling short of the optimal by at most 1%.

6. RELATED WORKS

Custom instruction selection for an application usually consists of two steps [28]. The initial step identifies a large set of candidate patterns from the program’s dataflow graph and their frequencies via profiling [4, 5, 10, 11, 18]. Given this library of patterns, the second step selects a subset to maximize the performance under different design constraints. Various approaches proposed for this step include dynamic programming [4], 0-1 Knapsack [11], greedy heuristic [9, 10], and ILP [22]. However, none of these approaches targets applications exploiting dynamic reconfiguration of custom functional units.

The major part of the research on temporal partitioning comes from the reconfigurable computing community. Usually, the partitioning is done at the task-level [6, 19, 8] while there could be some exceptions. Li et al. [23] partition at the loop level while Purna and Bhatia [26] perform partitions on the data flow graph. Moreover, some approaches do not consider software versions of the tasks. For example, Kaul et al. [19] propose a method in which a task graph is temporally partitioned, with the aim of minimiz-

ing overall latency. When directed acyclic task graphs are used as input, computing reconfiguration costs becomes simple. For example, Banerjee’s work [6] is able to reduce the partitioning problem as a scheduling problem because task graphs are used as input. In contrast, it is non-trivial to obtain the reconfiguration cost at the granularity of loops. It should be noted that while Purna and Bhatia’s work [26] partitions at the finer granularity of functions and operators, their work uses directed acyclic data flow graph as input as well.

Bondalapati and Prasanna [7] focus on mapping the statements within a loop into configurations to obtain a configuration sequence that gives the least execution time. While dynamic reconfiguration is used as well, their work focuses on intra-loop selection of configurations, i.e., their work operates on one loop only. Our work is different because not only do we consider multiple possible custom instructions set versions per loop, our algorithm allows for multiple loops within a configuration and some loops may remain in software. As such, our work is different from projects that explore the design space for individual loops such as [7].

Hardnett et al. [15] form a framework in which the dynamically reconfigurable architectural design space may be explored for specific applications. In particular, the register allocation problem is adapted to assign reconfigurable units to different custom instructions. An instruction scheduling algorithm for the custom instructions is implemented to minimize overall latency. While their architecture employs dynamically reconfigurable functional units, our work is differentiated from theirs in two specific areas. First, their custom instructions do not share the same functional unit, i.e., no spatial partitioning is required. Secondly, their work does not address the problem of reconfiguration cost directly. Rather, custom instructions are de-selected to relieve resource pressure rather than optimizing overall performance.

The work most related to our work is probably that of Li et al. [23] in which the Nimble compiler is implemented. Their work focuses on selecting loops from an application for hardware implementation while aiming to reduce dynamic reconfiguration overhead. Their work only considers a single loop in one configuration and they did not consider global reconfiguration cost when selecting loops to put in hardware.

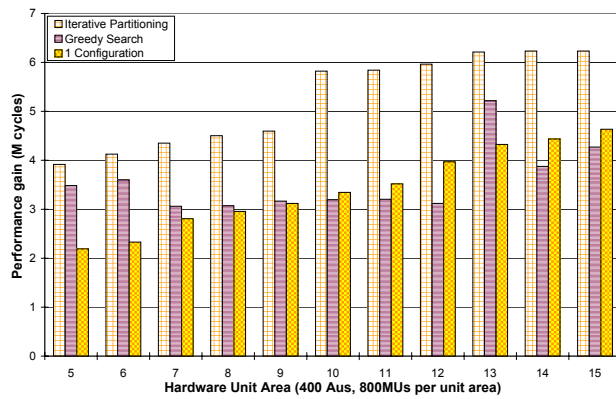
7. CONCLUSIONS

We have presented an algorithm to exploit dynamically configurable custom functional units for optimal performance gain. Given an input application, the algorithm selects and partitions the custom instructions corresponding to the loop kernels into different configurations that are reconfigured at run-time. The experimental results show that our algorithm is highly scalable while producing optimal or near-optimal performance gain.

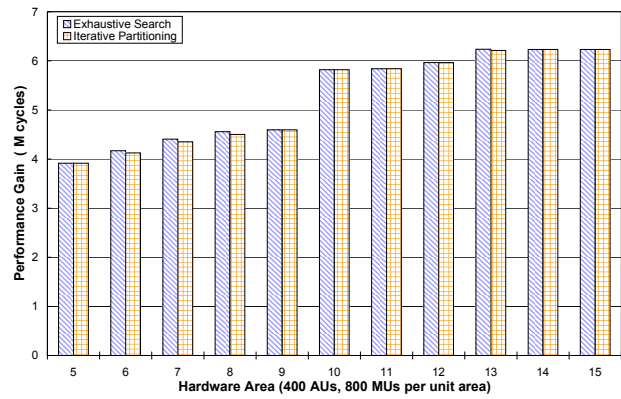
The work can be extended by considering configuration prefetching and partial reconfiguration. Previous work [24] have indicated that how early a configuration can be prefetched depends on several factors, including the relationship between the configurations and the placement of the prefetch instruction. Closely related to configuration prefetching is partial reconfiguration, which allows execution and reconfiguration of the fabric in parallel. In the future, we plan to extend our framework to handle these non-trivial issues.

8. ACKNOWLEDGMENTS

We thank the anonymous reviewers for their constructive comments and suggestions, which helped us to improve the paper. This work was partially supported by NUS research project R252-000-292-112 and A*Star SERC EHS-II project R-252-000-258-305.



(a) Comparison of Iterative partitioning, Greedy search, and the solution with no reconfiguration (1 configuration).



(b) Comparison of Exhaustive search and Iterative partitioning.

Figure 7: Experiment results for the case study of JPEG application.

9. REFERENCES

- [1] OpenIMPACT Compiler. <http://www.gelato.uiuc.edu/>.
- [2] ARC International. Customizing a soft microprocessor core.
- [3] J. M. Arnold. S5: The architecture and development flow of a software configurable processor. In *FPT*, 2005.
- [4] M. Arnold and H. Corporaal. Designing domain-specific processors. In *CODES*, 2001.
- [5] K. Atasu, L. Pozzi, and P. Inne. Automatic application-specific instruction-set extensions under microarchitectural constraints. In *DAC*, 2003.
- [6] S. Banerjee, E. Bozorgzadeh, and N. Dutt. Physically-aware HW-SW partitioning for reconfigurable architectures with partial dynamic reconfiguration. In *DAC*, 2005.
- [7] K. Bondalapati and V. K. Prasanna. Mapping loops onto reconfigurable architectures. In *FPL*, 1998.
- [8] K. S. Chatha and R. Vemuri. Hardware-software codesign for dynamically reconfigurable architectures. In *FPL*, 1999.
- [9] N. Cheung, S. Parameswaran, and J. Henkel. Inside: Instruction selection/identification & design exploration for extensible processors. In *ICCAD*, 2002.
- [10] N. Clark, H. Zhong, and S. Mahlke. Processor acceleration through automated instruction set customization. In *MICRO*, 2003.
- [11] J. Cong et al. Application-specific instruction generation for configurable processor architectures. In *FPGA*, 2004.
- [12] N. G. de Bruijn. *Asymptotic Methods in Analysis*. Dover Publications, 1981.
- [13] P. Faraboschi et al. Lx: A technology platform for customizable VLIW embedded processing. In *ISCA*, 2000.
- [14] R. E. Gonzalez. Xtensa: A configurable and extensible processor. *IEEE Micro*, 20(2), 2000.
- [15] C. Hardnett, K. V. Palem, and Y. Chobe. Compiler optimization of embedded applications for an adaptive soc architecture. In *CASES*, 2006.
- [16] G. Karypis and V. Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1998.
- [17] G. Karypis and V. Kumar. Multilevel k-way partitioning scheme for irregular graphs. *Journal of Parallel and Distributed Computing*, 48(1):96–129, 1998.
- [18] R. Kastner et al. Instruction generation for hybrid reconfigurable systems. *ACM TODAES*, 7(4), 2002.
- [19] M. Kaul, R. Vemuri, S. Govindarajan, and I. Ouass. An automated temporal partitioning and loop fission approach for FPGA based reconfigurable synthesis of DSP applications. In *DAC*, 1999.
- [20] B. W. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs. In *The Bell System Technical Journal*, volume 49(2), pages 291–307, 1970.
- [21] D. L. Kreher and D. R. Stinson. *Combinatorial Algorithms Generation, Enumeration and Search*. CRC Press Inc, 1998.
- [22] J. Lee, K. Choi, and N. Dutt. Efficient instruction encoding for automatic instruction set design of configurable asips. In *ICCAD*, 2002.
- [23] Y. Li, T. Callahan, E. Darnell, R. Harr, U. Kurkure, and J. Stockwood. Hardware-software co-design of embedded reconfigurable architectures. In *DAC*, 2000.
- [24] Z. Li and S. Hauck. Configuration prefetching techniques for partial reconfigurable coprocessor with relocation and defragmentation. In *FPGA*, 2002.
- [25] C. G. Nevill-Manning and I. H. Witten. Identifying hierarchical structure in sequences: A linear-time algorithm. *Journal Of Artificial Intelligence Research*, 7:67–82, 1997.
- [26] K. M. G. Purna and D. Bhatia. Temporal partitioning and scheduling data flow graphs for reconfigurable computers. *IEEE Transactions on Computers*, 48(6):579–590, 1999.
- [27] Stretch Inc. Stretch S5530 software configurable processor.
- [28] N. Topham. Challenges to automatic customization. In P. Inne and R. Leupers, editors, *Customizable Embedded Processors*. Morgan Kauffman, 2006.