

A C Compiler for a Processor with a Reconfigurable Functional Unit

Alex Ye, Nagaraj Shenoy, Scott Hauck, Prithviraj Banerjee
Department of Electrical and Computer Engineering, Northwestern University
Evanston, IL 60201, USA

Abstract

With the flexibility of the FPGA, reconfigurable systems are able to get significant speedups for some applications. Several architectures have been proposed to integrate a processor with an FPGA in the same chip [1][2][3][5]. The usage of the FPGA can be divided into two categories: FPGA as a coprocessor or FPGA as a functional unit [1][5]. The first scheme builds a large FPGA as a coprocessor [2][3]. The FPGA is usually equipped with some extra hardware to perform control flow operations and access memory. The tradeoff is less host processor area and larger communication overhead with the host processor. The second scheme builds a relatively simple FPGA as a reconfigurable functional unit (RFU). The function of the FPGA is less than the coprocessor scheme. On the other hand, the host processor gets more space and the FPGA gets a faster communication with the host.

This paper talks about a C compiler for Chimaera, which is a RISC/RFU architecture. It first describes three compilation techniques to extract instruction sequences from the applications to put into the RFU. It then shows some experimental results, which support both the compiler techniques and the architecture approach.

The first compilation technique is *Control Localization*. By changing some branches into a macroinstruction, several basic blocks can be combined into a large basic block. This increases the possibility of finding to find more sequences of instructions to put into the RFU. The second one is *SIMD Within A Register* (SWAR) optimization. The SWAR concept divides a register into several fields and performs several field-by-field operations in parallel. Most of the contemporary processors support a subset of SWAR in the form of multimedia extensions, such as Intel MMX, MIPS VIS, etc. The SWAR optimization in our compiler creates large instruction sequences to support a more general SWAR model. The third technique is *Instruction Combination*, which identifies all the possible instruction sequences to put into the RFU.

We have evaluated the system through a set of benchmarks. The results demonstrate that although the RFU doesn't provide functions for control flow operations and memory accesses, the compiler is still able to find enough instruction sequences to get significant speedups. The average speedup of all the benchmarks is 2.6. Sometimes the speedup is as high as 7.2. A notable fact is that most of the sequences are

small, which means that they are benefiting directly from the fast communication between the RFU and the host processor. This suggests that for some applications it is good to have a simpler FPGA with a faster communication. Furthermore, most of the short sequences are from complex address calculations. This indicates that this system has a wider application area than those of the typical reconfigurable systems, since many other applications have complex address calculations.

Compared to other Processor/FPGA works such as Garp[2], Napa [3], PRISC[5], etc, our system and PRISC are in the functional unit scheme category, while Garp and Napa belongs to the coprocessor scheme category. Garp and Napa target large portions of a program that can be implemented in the FPGAs. Usually a large portion will contain some instructions that cannot be implemented in the FPGA. As a result, the users may have to rewrite the program, or the compiler needs to find a way to execute the program without touching those instructions. Once such portions are found, the compiler has to schedule the computations both in time and space. While in our system and PRISC, the compiler's focus is on how to create and identify instruction sequences that can be changed into one RFU instruction. When the instructions are found, the compiler has to generate their mapping in the RFU. There is no need to schedule. Our work differs from PRISC mainly in the host processor part. The host processor in PRISC is a single-pipelined one. The major optimization of its compiler is a control flow conversion, which is not efficient in current architectures. The reason is that most of current architectures have branch prediction and can handle control flow operations better. In our system, we will consider the host to be a more advanced one. Therefore, our compiler is focused on the tasks that are not well handled by the advanced processors, e.g. instructions with sequential dependencies.

This paper makes three contributions. First, it describes two techniques that can effectively create more and larger instruction sequences for the RFU, namely Control Localization and SWAR. Second, it shows an algorithm to extract the instruction sequences in a basic block. Finally, it shows that the instruction sequences for the RFU exist in a wide range of applications and the Processor/RFU approach can achieve significant speedups.

1 Introduction

With the flexibility of the FPGA, reconfigurable systems are able to get significant speedups for some applications. As the general purpose processor and the FPGA each has its own suitable area of applications, several architectures are proposed to integrate a processor with an FPGA in the same chip. The usage of the FPGA in them can be divided into two categories: FPGA as a coprocessor or FPGA as a functional unit.

In the coprocessor schemes such as Garp[2] and Napa[3], the host processor is coupled with an FPGA based reconfigurable coprocessor. The coprocessor usually has the ability of accessing memory and performing control flow operations. There is a communication cost between the coprocessor and the host processor, which is several cycles or more. Therefore, these architectures tend to map a large portion of the application, e.g. a loop, into the FPGA. One calculation in the FPGA usually corresponds to a task that takes several hundred cycles or more.

In the functional unit schemes such as Chimaera[1] and PRISC[5], the host processor is integrated with an FPGA based Reconfigurable Functional Unit (RFU). One RFU Operation (RFUOP) can take on a task which usually requires several instructions on the host processor. As the functional unit is interfaced only with the register file, it cannot perform memory operations or control flow operations. The communication is faster than the coprocessor scheme. For example, in the Chimaera architecture, after an RFUOP's configuration is loaded, an invocation of it has no overhead in communication. This gives such architecture a larger range of application. Even in cases where only a few instructions can be combined into one RFUOP, we could still apply the optimization if the execution frequency is high enough.

In this paper, we talk about a C compiler for a RISC processor with an RFU. The target architecture is Chimaera. We describe how the compiler identifies sequences of statements in a C program and changes them into RFUOPs. We show the performance benefits that can be achieved by such optimizations over a set of benchmarks.

The paper is organized into five sections. In Section 2, we give an overview of the Chimaera architecture. Section 3 discusses the compiler organization and implementation in detail. In this section, we first discuss a technique to enhance the size of the instruction sequence: control localization. Next, we describe the application of the RFU to SIMD Within A Register (SWAR) operations. Lastly, we introduce an algorithm to identify RFUOPs in a basic block. Section 4 demonstrates some experimental results. We summarize this paper in Section 5.

2 Chimaera Architecture Overview

The overall Chimaera architecture is shown in Figure 1. The main component of the system is the Reconfigurable Functional Unit, which consists of FPGA-like logic designed to support high-performance computations. It gets inputs from the host processor's register file, or a shadow register file which duplicates a subset of the values in the host's register file.

The RFU contains several configurations at the same time. An RFUOP instruction will activate the corresponding configuration in the RFU. An RFU configuration itself determines from which registers it reads its operands. A single

RFUOP can read from all the registers connected to the RFU and then put the result on the result bus. The maximum number of input registers is 9 in Chimaera. Each RFUOP instruction is

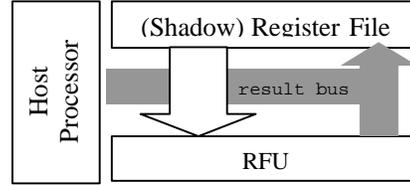


Figure 1. The overall Chimaera architecture

associated with a configuration and an ID. For example, an execution sequence “ $r2=r3<<2$; $r4=r2+r5$; $r6=lw\ 0(r4)$ ” can be optimized to “ $r4=RFUOP\ \#1$; $r6=lw\ 0(r4)$ ”. Here #1 is the ID of this RFUOP and “ $r5+r3<<2$ ” is the operation of the corresponding configuration. After an RFUOP instruction is fetched and decoded, the Chimaera processor checks the RFU for the configuration corresponding to the instruction ID. If the configuration is currently loaded in the RFU, the corresponding output is written to the destination register during the instruction writeback cycle. Otherwise, the processor stalls when the RFU loads the configuration.

3 Compiler Implementation

We have generated a C compiler for Chimaera which automatically maps some operations into RFUOPs. The

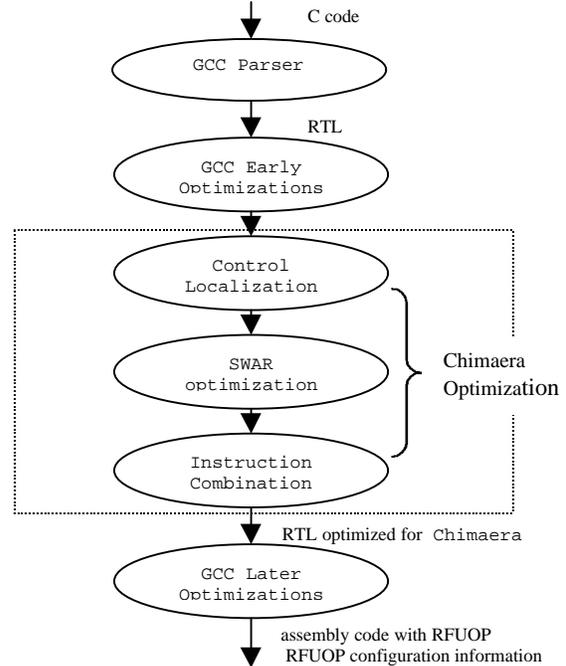


Figure 2: Phase ordering of the C compiler for Chimaera

generated code is run on a Chimaera simulator.

The compiler is built using the widely available GCC framework[6]. Figure 2 depicts the phase ordering of the implementation. The C code is parsed into the intermediate language of GCC: Register Transfer Language (RTL), which is then enhanced by several early optimizations such as common

expression elimination, flow analysis, etc. The partially optimized RTL is passed through the Chimaera optimization phase, as will be explained below. The Chimaera optimized RTL is then processed by later optimization phases such as instruction scheduling, registers allocation, etc. Finally, the code for the target architecture is generated along with RFUOP configuration information.

From the compiler’s perspective, we can consider an RFUOP as an operation with multiple register inputs and a single register output. The goal of the compiler is to identify the suitable multiple-input-single-output sequences in the programs and change them into RFUOPs.

Chimaera Optimization consists of three steps: Control Localization, SWAR optimization and Instruction Combination. Due to the configuration loading time, these optimizations can be applied only in the kernels of the programs. Currently, we only optimize the innermost loop in the programs.

The first step of Chimaera optimization is control localization. It will transform some branches into one macroinstruction to form a larger basic block. The second step is the *SIMD Within A Register* (SWAR) Optimization. This step searches the loop body for subword operations and unrolls the loop when appropriate. The third step is instruction combination. It takes a basic block as input and extracts the multiple-input-single-output patterns from the data flow graph. These patterns are changed into RFUOPs if they can be implemented in RFU. The following subsections discuss the three steps in detail.

3.1 Control Localization

In order to get more speedup, we want to find larger and more RFUOPs. Intuitively, a larger basic block contains more instructions, thus has more chances of finding larger and more RFUOPs. We find that control localization technique [4][5] is useful in increasing the size of basic blocks. Figure 3 shows an example of it. After control localization, several branches are combined into one macroinstruction, with multiple output and multiple input. In addition to enlarging the basic block, the control localization sometimes finds RFUOPs directly. When a macroinstruction has only one output, and all the operations in it can be implemented in the RFU, this macroinstruction can be mapped into an RFUOP. This RFUOP can speculatively compute all operations on different branch paths. The result on the correct path where the condition evaluates to true is selected to put into the result bus. This macro instruction is called as “CI macroin” and can be optimized by Instruction Combinaton.

3.2 SWAR Optimization

As a method to exploit medium-grain data parallelism, SIMD (single instruction, multiple data) has been used in parallel computers for many years. Extending this idea to general purpose processors has led to a new version of SIMD, namely *SIMD Within A Register* (SWAR)[7]. The SWAR model partitions each register into fields that can be operated on in parallel. The ALUs are set up to perform multiple field-by-field operations. SWAR has been successful in improving the multimedia performance. Most of the implementations of this concept are called multimedia extensions, such as Intel

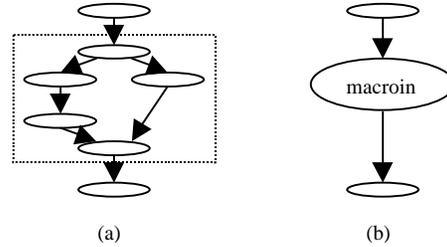


Figure 3:

- (a) control flow graph before control localization. Each oval is an instruction, and the dashed box marks the code sequence to be control localized.
- (b) control flow graph after control localization

MMX, HP MAX, SUN SPARC VIS, etc. For example, “PADDB A, B” is an instruction from Intel MMX. Both operands A and B are 64-bit and are divided into eight 8-bit fields. The instruction performs eight additions in parallel and stores the eight results to A.

However, current implementations of SWAR do not support a general SWAR model. Some of their limitations are:

- The input data must be packed and aligned correctly, causing packing and unpacking penalties sometimes.
- Most of current hardware implementations support 8, 16 and 32-bit field size only. Other important sizes such as 2-bit and 10-bit are not supported.
- Only a few operations are supported. When the operation for one item becomes complex, SIMD is impossible. For example, the following code does not map well to a simple sequence of SIMD operations:

```
char out[100], in1[100], in2[100];
for(i=0; i<100; i++) {
    if ((in1[i]-in2[i])>10)
        out[i]=in1[i]-in2[i];
    else
        out[i]=10;}

```

With the flexibility of the FPGA, the RFU can support a more general SWAR model without the above disadvantages. The only requirement is that the output fields should fit within a single register. The inputs don’t need to be stored in packed format, nor is there limitation on the alignment. In addition, complex operations can be performed. For example, the former example can be implemented in one RFUOP.

Our compiler currently supports 8-bit field size, which is the size of “char” in C. In current implementation, the compiler looks for the opportunity to pack several 8-bit outputs into a word. In most cases, this kind of pattern exists in the loop with stride one. Therefore, the compiler searches for the pattern such that the memory store size is a byte and the address changes by one for each loop iteration. When such patterns are found, the loop is unrolled four times. In the loop unrolling, conventional optimizations such as local register renaming and strength reduction are performed. In addition, the four memory stores are changed to four sub-register movements. For example,

```
“ store_byte r1, address;
  store_byte r2, address+1;
  store_byte r3, address+2;
  store_byte r4, address+3;”
are changed into
“(r5, 0)=r1; (r5, 1)=r2;
 (r5, 2)=r3; (r5, 3)=r4;”.
```

The notation (r, n) refers to the n th byte of register r . We generate a pseudo instruction "collective-move" that moves the four sub-registers into a word register, e.g. " $r5=(r5,0)(r5,1)(r5,2)(r5,3)$ ". In the data flow graph, the four outputs merge through this "collective-move" into one. Thus a multiple-input-single-output subgraph is formed. The next step, Instruction Combination, can recognize this subgraph and change it to an RFUOP when appropriate. Finally, a memory store instruction is generated to store the word register. The compiler then passes the unrolled copy to the instruction combination step.

```

ALGORITHM: FindSequences
INPUT:DFG G, NCI, CI, Live-on-exit registers R
OUTPUT: A set of RFU sequences S
begin
  S= $\emptyset$ 
  for each node  $n \in G$ 
    Color( $n$ )  $\leftarrow$  WHITE
  end
  for each node  $n \in NCI$ 
    Color( $n$ )  $\leftarrow$  BLACK
    for each node  $p \in \text{Pred}(n)$ 
      if  $p \in CI$  then
        Color( $p$ )  $\leftarrow$  BROWN
      endif
    end
  end
  for each register  $r \in R$ 
     $n \leftarrow$  the last node that updates  $r$  in  $G$ 
    if  $n \in CI$  then
      Color( $n$ )  $\leftarrow$  BROWN
    endif
  end
  for each node  $n \in G$ 
    if Color( $n$ )=BROWN then
      sequence= $\emptyset$ 
      AddSequence( $n$ , sequence)
      if sizeof(sequence)>1 then
        S= $S \cup$ {sequence}
      endif
    endif
  end
end

AddSequence( $n$ , sequence)
begin
  if Color( $n$ )=(BROWN or WHITE) then
    sequence $\leftarrow$ sequence $\cup$ { $n$ }
    for each  $p \in \text{Pred}(n)$ 
      AddSequence( $p$ , sequence)
    end
  endif
end

```

Figure 3: Algorithm to find RFU sequences

3.3 Instruction Combination

The instruction combination step analyzes a basic block and changes the RFU sequences into RFUOPs. It first finds out what instructions can be implemented in the RFU. It then identifies the RFU sequences. At last, it selects the appropriate RFU sequences and changes them into RFUOPs.

We categorize instructions into Chimaera Instruction (CI) and Non-Chimaera Instruction (NCI). Currently CI includes logic operation, constant shift and integer add/subtract. The "collective_move", "subregister movement" and "CI macroin" are also considered as CI. NCI includes other instructions such as multiplication/division, memory load/store, floating-point operation, etc.

The algorithm FindSequences in Figure 3 finds all the maximum instruction sequences for the RFU. It colors each node in the data flow graph(DFG). The NCI instructions are marked as BLACK. A CI instruction is marked as BROWN when its output must be put into a register, that is, the output is live-on-exit or is the input of some NCI instructions. Other CI instructions are marked as WHITE. The RFU sequences are the subgraphs in the DFG that consists of BROWN nodes and WHITE nodes.

The compiler currently changes all the identified sequences into RFUOPs. Under the assumption that every RFUOP takes one cycle and the configuration loading time can be amortized over several executions, this gives an upper bound of the speedup we could expect from Chimaera. In the future, we will take into account other factors such as the FPGA size, configuration loading time, actual RFUOP execution time, etc.

4 Experimental Results

We have tested the compiler's output through a set of benchmarks on the Chimaera simulator. The simulator is a modification of SimpleScalar Simulator[6]. The simulated architecture has 32 general purpose 32-bit registers and 32 floating point registers. The instruction set is a superset of MIPS-IV ISA. Presently, the simulator executes the programs sequentially and gathers the instruction count. Our future goal is to obtain more realistic results, considering multiple pipelines, configuration loading time and out-of-order execution.

Early results on some benchmarks are presented in this section. Each benchmark is compiled in two ways: one is using "gcc -O2", the other is using our Chimaera compiler. We studied the differences between the two versions of assembly codes as well as the simulation results. In the benchmarks, decompress.c and compress.c are from Honeywell benchmark[10], jacobi and life are from Raw benchmark[9], image reconstruction[8] and dct[11] are implementations of two program kernels of MPEG, image restoration is an image processing program. They are noted as dcmp, cmp, life, jcbi, dct, rcn and rst in the following figure.

	insn1	insn2	Speedup	IC	CL	SWAR
dcmp	37.4M	32.7M	1.14	100%	0	0
cmp	34.2M	28.6M	1.20	80%	20%	0
jcbi	2.1M	1.3M	1.63	94%	6%	0
life	6.7M	5.0M	1.34	100%	0	0
dct	78.4M	52.6M	1.49	10%	82%	8%
rcn	11M	2.6M	4.23	6%	10%	84%
rst	22.3M	3.1M	7.19	16%	0	84%
avg	-	-	2.60	22%	7%	61%

Figure 4: Performance results over some benchmarks. The "avg" row is the average of all benchmarks.

Figure 4 shows the simulation results of the RFU optimizations. Insn1 and insn2 are the instruction counts without and with RFU optimization. The speedup is calculated as

insn1/insn2. The following three columns IC, CL and SWAR stand for the portion of performance gain from Instruction Combination, Control Localization and SWAR respectively. The three optimizations give an average speedup of 2.60. The best speedup is up to 7.19.

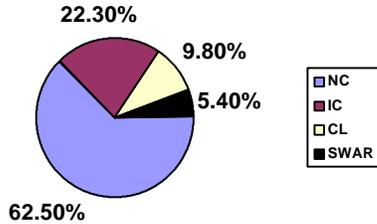


Figure 5: Distribution of the kernel instructions

To illustrate the impact of each optimization on the kernel sizes, we categorize instructions into four types: *NC*, *IC*, *CL* and *SWAR*. *NC* is the part of instructions that cannot be optimized for Chimaera. *NCI* instructions and some non-combinable integer operations fall in this category. *IC*, *CL* and *SWAR* stand for the instructions that can be optimized by Instruction Combination, Control Localization and SWAR optimization respectively. Figure 5 shows the distribution of these four types of instructions in the program kernels. After the three optimizations, the kernel size can be reduced by an average of 37.5%. Of this amount, 22.3% is from Instruction Combination, 9.8% from Control Localization and 5.4% from SWAR.

Further analysis shows that 58.4% of the *IC* portion comes from address calculation. For example, the following C code "int a[10], ...=a[i]" is translated to "r3=r2<<2, r4=r3+r1, r5=lw 0(r3)" in assembly. The first two instructions can be combined in Chimaera. The large portion of address calculation indicates that our optimizations can be applied to a wide range of applications, as long as they have complex address calculations in the kernel. Furthermore, as the address calculation is basically sequential, existing ILP architectures like superscalar and VLIW cannot take advantage of it. This suggests that we may expect speedup if we integrate a RFU into an advanced ILP architecture.

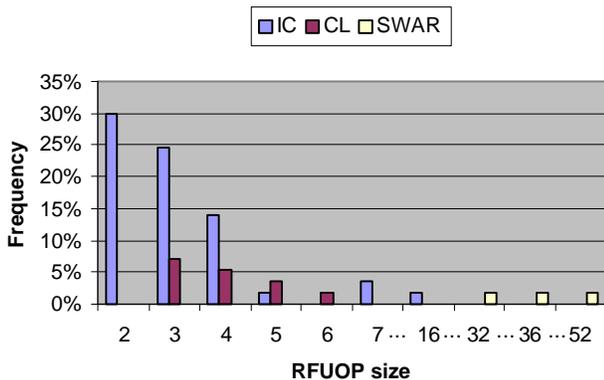


Figure 6: The frequency of RFUOP sizes. IC, CL and SWAR are the RFUOPs generate by Instruction Combination, Control Localization and SWAR respectively.

Figure 6 illustrates the frequencies of different RFUOP sizes. For Instruction Combination and Control Localization,

most of the sizes are from 2 to 6. These small sizes indicate that these techniques are benefiting from the fast communication of the functional unit scheme. In the coprocessor scheme, the communication overhead would make them prohibitive to apply.

In summary, the results show that the compilation techniques are able to create and find many instruction sequences for the RFU. Most of their sizes are several instructions, which demonstrate that the fast communication is necessary. The system gives an average speedup of 2.6.

5 Conclusion

This paper describes a C compiler for the Processor/FPGA architecture when the FPGA is served as a Reconfigurable Functional Unit (RFU).

We have introduced an instruction combination algorithm to identify RFU sequences of instructions in a basic block. We have also shown that the control localization technique can effectively enlarge the size of the basic blocks and find some more sequences. In addition, we have illustrated the RFU support for SWAR. By introducing "sub-register movement" and "collective-move", the instruction combination algorithm is able to identify complex SIMD instructions for the RFU.

Finally, we have presented the experimental results which demonstrate that these techniques can effectively create and identify larger and more RFU sequences. With the fast communication between RFU and the processor, the system can achieve considerable speedups.

References

- [1] Omitted for anonymous reviewing
- [2] John R. Hauser and John Wawrzynek. GARP: A MIPS processor with a reconfigurable coprocessor. Proceedings of IEEE Workshop on FPGAs for Custom Computing Machines (FCCM), Napa, CA, April 1997.
- [3] M.B. Gokhale, et al. Napa C: Compiling for a Hybrid RISC/FPGA Architecture, FCCM 98, CA, USA
- [4] Walter Lee, Rajeev Barua and et al. Space-Time Scheduling of Instruction-Level Parallelism on a Raw Machine, MIT. ASPLOS VIII 10/98, CA, USA
- [5] R. Razdan, PRISC: Programmable Reduced Instruction Set Computers, Ph.D. Thesis, Harvard University, Division of Applied Sciences, 1994
- [6] <http://www.cs.wisc.edu/~mscalar/simplescalar.html>
- [7] Paolo Faraboschi, et al. The Latest Word in Digital and Media Processing, IEEE signal processing magazine, Mar 1998.
- [8] Selliah Rathnam, et al. Processing the New World of Interactive Media, IEEE signal processing magazine March 1998
- [9] J.Babb, M.Frank, et al. The RAW benchmark Suite: Computation Structures for General Purpose Computing. FCCM, Napa Vally, CA, Apr.1997
- [10] Honeywell Inc, <http://www.htc.honeywell.com/projects/acsbench/>
- [11] Koichi Atsuta, ET5Lab Compression Group, <http://marine.et.u-tokai.ac.jp/database/koichi.html>.
- [12] Compiling For SIMD Within A Register, Randall J. Fisher and Henry G. Dietz, 1998 Workshop on Languages and Compilers for Parallel Computing, North Carolina, Aug 1998