

Merged Dictionary Code Compression for FPGA Implementation of Custom Microcoded PEs

BITA GORJIARA, MEHRDAD RESHADI, and DANIEL GAJSKI

University of California, Irvine

Horizontal Microcoded Architecture (HMA) is a paradigm for designing programmable high-performance processing elements (PEs). However, it suffers from large code size, which can be addressed by compression. In this article, we study the code size of one of the new HMA-based technologies called No-Instruction-Set Computer (NISC). We show that NISC code size can be several times larger than a typical RISC processor, and we propose several low-overhead dictionary-based code compression techniques to reduce its code size. Our compression algorithm leverages the knowledge of “don’t care” values in the control words and can reduce the code size by 3.3 times, on average. Despite such good results, as shown in this article, these compression techniques lead to poor FPGA implementations because they require many on-chip RAMs. To address this issue, we introduce an FPGA-aware dictionary-based technique that uses the dual-port feature of on-chip RAMs to reduce the number of utilized block RAMs by half. Additionally, we propose cascading two-levels of dictionaries for code size and block RAM reduction of large programs. For an MP3 application, a merged, cascaded, three-dictionary implementation reduces the number of utilized block RAMs by 4.3 times (76%) compared to a NISC without compression. This corresponds to 20% additional savings over the best single level dictionary-based compression.

Categories and Subject Descriptors: C.3 [Special-Purpose and Application-Based Systems] Microprocessor/microcomputer applications

General Terms: Algorithms, Design, Performance, Experimentation

Additional Key Words and Phrases: Microcoded architectures, no-instruction-set computer, memory optimization, dictionary based compression, FPGA

ACM Reference Format:

Gorjiara, B., Reshadi, M., and Gajski, D. 2008. Merged dictionary code compression for FPGA implementation of custom microcoded PEs. *ACM Trans. Reconfig. Techn. Syst.* 1, 2, Article 11 (June 2008), 21 pages. DOI = 10.1145/1371579.1371583. <http://doi.acm.org/10.1145/1371579.1371583>.

Authors’ address: Center for Embedded Computer Systems, University of California at Irvine; email: {gorjiara, reshadi, gajski}@cecs.uci.edu.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permission may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701, USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2008 ACM 1936-7406/2008/06-ART11 \$5.00 DOI: 10.1145/1371579.1371583. <http://doi.acm.org/10.1145/1371579.1371583>.

ACM Transactions on Reconfigurable Technology and Systems, Vol. 1, No. 2, Article 11, Pub. date: June 2008.

1. INTRODUCTION

Shrinking time-to-market and high demand for productivity has driven traditional hardware designers to use design methodologies that start from high-level languages. However, meeting design constraints of automatically generated Processing Elements (PEs) is often a challenging and time-consuming task for designers. Moreover, slight changes in the high-level specification require rerunning the behavioral synthesis tools, producing a new datapath, and redoing the physical design process. To avoid repeating timing closure and physical synthesis phases, a new generation of custom-PE design technology that is capable of both generating custom datapaths as well as reprogramming existing ones (without further modifications) is developed. In these technologies, first a custom datapath is generated for an application, and then the datapath is synthesized and laid out properly to meet timing and physical constraints. The final step is to compile the program on the generated datapath. If the application is changed after synthesis, it is simply recompiled on the existing datapath. This feature significantly improves the productivity of the designer by preventing repetition of physical synthesis phase. Examples of such technology include ARM OptimoDE, No-Instruction-Set Computer (NISC) [Reshadi and Gajski 2005; Reshadi et al. 2005], and TIPI [Weber and Keutzer 2005]. These techniques are targeted for statically scheduled Horizontal Microcoded Architectures (HMA) [Agrawala and Rauscher 1976].

A microcode is a set of bits that controls the units of datapath for one cycle. In statically scheduled HMAs, the compiler compiles the program directly to microcode without using instruction abstraction. HMAs can potentially have better performance, lower power, and lower area than conventional instruction-based processors. This is due to giving the compiler more fine-grained control over the datapath, and hence utilizing datapath resources more efficiently. As a result, highly parallel architectures can be designed as HMA without any concern about the complexity of the controller, hardware scheduler, and instruction decoder. Despite all these benefits, HMAs suffer from “code bloating.”

This article studies FPGA-implementation of soft-core HMA-based PEs. We compare the code size of a new HMA-based design methodology, called NISC, with that of traditional RISC processors. We observed that although NISC PEs outperform typical RISC processors by five times on average, their code sizes are about four times larger than those of RISC. In this article, we propose several low-overhead dictionary-based code compression techniques to reduce the code size. Our compression algorithm leverages the knowledge of “don’t care” values in the control words to improve the compression efficiency and can reduce the code size by 3.3 times, on average. Despite such good results, as shown in this article, these compression techniques lead to poor FPGA implementations because they require many on-chip RAMs. To overcome this limitation, we propose to merge every two dictionaries into a single dual-port memory unit on FPGAs. Using this approach, the block RAM utilization is improved by 46%. Also, for large applications, we propose using *cascaded dictionaries*, where multi-levels of dictionaries are used to decompress the code. For MP3

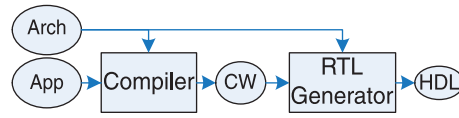


Fig. 1. Flow of our toolset.

application, a merged, cascaded, three-dictionary implementation reduces the number of utilized block RAMs by 4.3 times (76%) compared to a NISC without compression. This corresponds to 20% additional savings compared to the best single-level dictionary-based compression.

This article is organized as follows: Section 2 presents an overview of NISC Technology. Section 3 presents a motivating example to emphasize the need for code-size reduction techniques in HMAs. Section 4 is an overview of existing code-size reduction techniques. Section 5 discusses our multi-dictionary compression approach, its effectiveness in terms of compression ratio, and its limitation in terms of number of utilized block RAMs. Section 6 introduces our FPGA-aware compression technique that can significantly improve block RAM utilization in FPGAs. Section 7 proposes cascading two-level of dictionaries for code optimization of large program. Section 8 discusses the worst-case performance penalty of decompression. Finally, Section 9 concludes the article.

2. OVERVIEW OF NISC TECHNOLOGY

In the NISC design flow [Reshadi et al. 2005; Gorjiara et al. 2006], a custom architecture is generated or selected for a given application and then, the program is compiled on the architecture to generate low-level microcodes that we call control words (CW). Finally, the HDL code of the NISC is generated in register-transfer level (RTL) according to the architecture description and the output control words [Gorjiara et al. 2006]. Our toolset (Figure 1) is available online at <http://www.cecs.uci.edu/~nisc>, where users can specify a new NISC architecture using our Architecture Description Language (ADL) called Generic Netlist Representation (GNR) [Gorjiara et al. 2006] and compile their application on it. Also, automatic tools can be used to generate NISC architectures customized for one or more applications [Gorjiara and Gajski 2008; Gorjiara 2007; Trajkovic et al. 2006]. NISC customization can be done in three ways: (1) by changing number and type of components and their interconnectivity; [Gorjiara and Gajski 2008] (2) by adding custom functional units (ranging from simple bitwise operations to complex multi-operand operations); [Reshadi et al. 2008] (3) by adding one or more external accelerators [Reshadi and Gajski 2007; Reshadi 2007], which may be custom NISCs themselves. The difference between (2) and (3) is that in (2) compiler schedules the communication between the custom unit and the rest of the datapath, while in (3) the software should explicitly define the communication. Customization can significantly improve performance and code size of a design [Gorjiara 2007], but it is out of scope of this article. In this article, we focus only on compression-based code size reduction techniques.

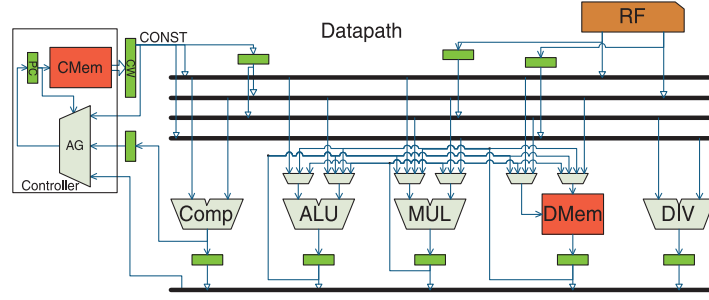


Fig. 2. Block diagram of GNISC architecture.

RF_ra0	RF_ra1	RF_wa	RF_we	ALU_op	Mux0_sel	Mux1_sel	R0_load	...	Constant
--------	--------	-------	-------	--------	----------	----------	---------	-----	----------

Fig. 3. Fields of a control word.

Figure 2 shows an example of a NISC PE, which consists of a datapath and a controller. The datapath contains functional units, register file, registers, multiplexers, and memory. This datapath is general enough to run many applications; therefore, we call it General NISC (GNISC). Our approach relies on a sophisticated compiler [Reshadi 2007] to compile a program described in a high-level language to control words that directly drive the control signals of components in the datapath. The CWs are stored in a control memory (CMem) in programmable PEs, or they are synthesized to lookup-table logic in hardwired dedicated PEs. Figure 3 shows the fields of a sample control word. Corresponding to each control signal of each component in the datapath, a field is added to the control words. Also, one or more constant fields are added to store constants in the program. Our compiler generates “0,” “1,” or “don’t care” values for the bits of the control words. A “don’t care” value (denoted by “X”) indicates that the corresponding unit is idle at a given cycle and its control signal can be assigned to “0” or “1” without affecting program behavior.

Compared to traditional RISC processors, NISC architecture does not have an instruction decode stage. Also, the operations are statically scheduled; hence, it does not have a hardware scheduler. As a result, highly parallel architectures can be designed using NISC without any concern about the complexity of instruction decoder and hardware scheduler.

Similar to NISC, VLIW/EPIC processors [Rau et al. 1989; Codwell et al. 1987] are also statically scheduled. However, they have several differences:

- (1) VLIWs have instruction decoders because they still have the concept of instructions. In fact their decoders are larger than those of RISC processors because they should decode several RISC instructions simultaneously. For example, in TI TMS320C62xx (a modern VLIW architecture), the instruction decoder consists of two pipeline stages. In contrast, NISC does not need any instruction decode stage.

Table I. Area and Clock Frequency of MicroBlaze and GNISC

Processors	Clock freq.(MHz)	# 4-input LUTs
MicroBlaze	105	1581
GNISC	100	1576

- (2) Adding a new unit to VLIW requires adding a slot to its wide instructions (each slot is equivalent of a RISC instruction). Also the decoder unit must be extended accordingly. However, in NISC, for each new unit, only one or a few control fields are added to the control word. Therefore, more units can be added to the datapath at little cost.
- (3) Resource sharing can be implemented more efficiently in NISC than VLIWs, because NISC compiler has low-level control over elements of the datapath. For example, instead of adding two constant fields to the control word of GNISC architecture in Figure 2, one constant field and an *assistant register* is added. If in a given cycle, two constants are needed in the datapath, the compiler can schedule one of the constants a few cycles earlier and transfer it to the assistant register for future use. Similarly, assistant registers are added at the output of register file to reduce the number of required register-file read ports. Such optimizations are to some extent similar to renaming and reservation stations in superscalar processors; but in NISC, they are done at compile-time and do not impose any hardware overhead. In contrast, these optimizations are very challenging in VLIWs because many low-level instructions must be designed for transferring data to the internal registers.

Despite their differences, both NISC and VLIWs have code size issue. The code size reduction techniques proposed for VLIW processors are applicable to NISC as well. Before discussing these techniques, we first present a motivational example that compares the quality and code size of a NISC architecture with those of an instruction-based processor. Since we did not have access to toolset and HDL description of a VLIW processor, we compared NISC to a well-known RISC processor.

3. MOTIVATIONAL EXAMPLE

In this section, we compare the performance and code size of the GNISC with a similar-size RISC processor. Since our toolset generates synthesizable code for Xilinx FPGAs, we choose the Xilinx MicroBlaze for comparing an instruction-based processor with a microcoded one. We synthesized both processors on a Xilinx Virtex4 (90nm) FPGA package using ISE 8.2. We configured MicroBlaze to have an integer multiplier and a divider (no barrel shifter, no floating point unit). Table I shows the area (in terms of 4-input LUTs) and clock frequency of the processors. Both processors run at about 100MHz, and occupy nearly the same number of 4-input LUTs.

We compiled and simulated a set of benchmarks including *dijkstra*, *sha*, *adpcm_coder*, *adpcm_decoder* and *CRC32* from MiBench (the free version of EEMBC embedded benchmarks at <http://www.eecs.umich.edu/mibench>), and a fixed-point Mp3 decoder (more than 10,000 lines of C code). For all benchmark,

Table II. Comparing GNISC with MicroBlaze

Benchmarks	MicroBlaze			GNISC			GNISC MicroBlaze	
	#cycles	code size		#cycles	code size		speedup (x)	code size ratio
		KB	#BRAM		KB	#BRAM		
adpcm_coder	256748693	1.956	1	74321930	6.960	4	3.45	5.10
adpcm_decoder	322766405	1.364	1	63082673	5.075	3	5.12	2.59
CRC32	209436647	1.264	1	21901993	2.567	3	9.56	2.03
dijkstra	25927532	1.928	1	9764682	9.614	6	2.66	4.99
sha	183030479	3.156	2	19282976	14.123	11	9.49	4.47
Mp3	2668445	44.62	21	897452	216.659	117	2.97	4.86
Average							5.54	4.01

we have removed file I/O and *printf* calls to make the code suitable for FPGAs. MiBench provides a small and a large input for the benchmarks. We used a subset of small input. We also set the compiler optimizations to the maximum possible level (i.e., -O2) to achieve the best performance with both NISC and MicroBlaze. For each benchmark, to get the accurate execution cycle count, we generated RTL Verilog code of the design and simulated it using Modelsim simulator.

Table II shows the number of cycles and code size of each benchmark on the two processors. The code size of MicroBlaze (the third column) is the size of instruction section (.text) of the ELF file generated by the compiler. Also the fourth column is the code size in terms of number of on-chip Block RAMs (BRAMs). BRAMs are ASIC memory units that exist on modern FPGA chips. Similarly, the sixth and seventh columns show the code size of GNISC in terms of KB and number of utilized BRAMs. The eighth column shows speedup of GNISC compared to MicroBlaze. The ninth column shows the ratio of GNISC code size (KB) to that of MicroBlaze. On average, GNISC runs 5.54 times faster than MicroBlaze, while its code size is four times larger. Although GNISC occupies almost the same number of LUTs as of MicroBlaze, it needs significantly more BRAMs to store the code. In this article, we show that different dictionary-based compression techniques can reduce the code size (in terms of KB), but they may fail to reduce the number of utilized BRAMs in FPGA-based implementation. The goal of our code optimization technique is to reduce the code size of NISC processors (both in terms of BRAMs and KB) while maintaining their performance benefits.

4. OVERVIEW OF EXISTING CODE-SIZE REDUCTION TECHNIQUES

In general-purpose processors, the instruction-set abstraction is used to reduce the code size of processors. In RISC processors, designers define 32-bit or 16-bit [Segars et al. 1995; Grehan 1997] instructions to encode wide control words. At runtime, the instructions are decoded back to the control words using a hardware decoder. In most processors, one or more pipeline stages are added to the datapath for instruction decoding. As a result, instructions increase branch delay, and hence, affect the performance of the processor (branch delay is the number of cycles between the fetch of a branch instruction and finishing the computation of branch target address). Branch prediction can partially address this issue, but it increases the area of the processor. On the other

hand, designing an instruction set is a very complex and time-consuming task for a typical hardware designer; because the compiler, assembler, linker and instruction decoder must be redesigned to handle custom instructions. To increase the productivity of designers and to give more control to the compiler, in our approach, we eliminate the need for instruction-set design, and compile the application directly to control words. To reduce the code size, instead of using instructions, we directly compress the control words. This way, we replace instruction decode stage(s) with control-words decompression stage(s). In other words, the performance penalty of decompression is the same as the instruction decoder.

As discussed in Section 2, VLIWs have also very large code sizes. In traditional (aka canonical) VLIWs, instructions have several slots, each of which corresponds to an execution unit in the datapath. Due to limited parallelism in the application, at every cycle, some of the units are idle; hence, their corresponding instruction slot is filled with NOP operations. The existence of many NOPs in the code increases its size significantly. To address this issue, modern commercially successful VLIW architectures (such as TMS320C6x) allow more flexible instructions based on the idea of Various Length Execution Set (VLES). In this approach, NOPs are removed as much as possible and a few shortened instructions are packed into one wide fetch packet. Different instruction packing algorithms are discussed by Saghir [1998] and Wang [2001]. For unpacking VLES instructions, a few pipeline stages are added to the processor. For example, in TMS320C6x, fetch/unpacking consists of four pipeline stages; hence, in this processor, fetch and decode need overall six pipeline stages. If compression is applied to VLIWs, it may also need a few more pipeline stages for decompression. Such high number of pipeline stages increases branch delay and area significantly. In contrast, NISC does not have instructions or instruction slot. Therefore, it does not need the instruction decode stages. Also, since it does not have NOP instructions, it does not need packing, and hence unpacking pipeline stages. The only thing it needs is compression, which may add a few pipeline stages to the processor for decompression.

Compression algorithms [Rafail 1994] can be categorized to two main groups: dictionary-based (DBC) and arithmetic (statistic) based (ABC). Also, they can be categorized as fixed length or variable length depending on whether the compressed words (aka codewords) have the same length or not.

The dictionary-based compression techniques rely on the fact that the same patterns appear many times in the data. Figure 4 shows how these techniques typically work. The unique patterns are stored in a dictionary and the original data is replaced with indexes to the dictionary. If original data has N words and each word is n bits, then $N \times n$ bits must be stored. If the number of unique patterns is M , and $m = \log_2 M$ is much smaller than n , then, a dictionary based technique can compress the data down to $N \times m + M \times n$ bits, which is often smaller than $N \times n$. However, the decompression costs one extra lookup. To improve the compression, the dictionary-based compression can be combined with Huffman coding where the frequently repeating instructions are placed in low addresses of the dictionary and are coded with fewer number of bits. In this case, the codewords become variable length and need more cycles to

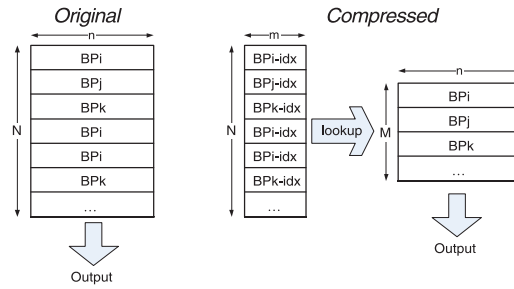


Fig. 4. Dictionary-based code compression.

decompress. Different variations of dictionary-based compression are proposed for RISC processors. The Code Compressed RISC Processor (CCRP) [Wolfe and Chanin 1992] combines dictionary-based compression with Huffman coding. IBM's CodePack [Kemp et al. 1998; Lefurgy et al. 1999] improves the compression efficiency of CCRP by partitioning instructions to two halves and using two dictionaries to store the unique patterns of each half. Corliss et al. [2003], Fraser [2002], and Lau et al. [2003] extend the concept of dictionary-based compression to sequence of instructions. In these approaches, unique *sequences* of instructions are identified and stored in a dictionary. Dictionary-based compression has also been applied to VLIW processors. Ishiura and Yamaguchi [1997] extend the instruction partitioning approach in CodePack by automatically partitioning instructions to several fields so that the overall code size is minimized. This approach is applied to traditional (canonical) VLIW processors and shows two to three times reduction in code size. However, in modern (VLES-style) VLIW processors dictionary-based compression is very ineffective (only 10–15% code reduction) because instruction packing increases the number of dictionary entries. To improve its efficiency, Ros and Sutton [2004] propose combining nearly-identical instructions (that are different only in a few bits) into a single entry in the dictionary. However, they add a new field to codewords to specify the bits that must be toggled during decompression. This approach improves the compression efficiency by a few percentages. Since dictionary-based compression is less effective for modern VLIW processors, arithmetic-based compression techniques (ABC) are proposed [Xie et al. 2001; 2002], which have better compression ratio but have significantly higher decompression overhead.

Depending on the decompression overhead, compression algorithms can be implemented in two ways for general-purpose processors [Xie et al. 2002]:

- (1) Pre-cache: where the code is decompressed between main memory and cache. In other words, main memory contains the compressed code, and the cache contains the decompressed code. In this approach, the decompression penalty is paid only when a cache miss occurs. Compression techniques that have relatively high decompression overhead (more than one or two cycles) should be implemented as pre-cache. Huffman-coded dictionary based compression techniques [Wolfe and Chanin 1992; Kemp et al.

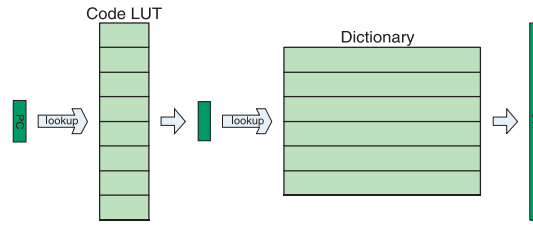


Fig. 5. One-dictionary code compression (*opt1*).

1998; Lefurgy et al. 1999; Corliss et al. 2003; Fraser 2002; Prakash et al. 2003] and most arithmetic-based techniques [Xie et al. 2001; 2002] fall into this category.

- (2) Post-cache: where the code is decompressed between cache and processor. In other words, both cache and main memory contain compressed code. In this approach, the decompression hardware is on the critical path and should be very fast. Some of the dictionary-based technique [Ishiura and Yamaguchi 1997; Ros and Sutton 2004; Xie et al. 2003] and one of the arithmetic based techniques [Xie et al. 2002] fall into this category.

In custom reprogrammable hardware units, the entire program and data may fit in on-chip memory blocks. Therefore, cache may impose unnecessary overhead. In our approach, we limit the decompression overhead to one or two cycles in order to allow pre-cache, post-cache, and no-cache implementations. Compared to previous approaches, our approach has several differences:

- (1) We leverage the existence of “don’t care” values in our binary to improve efficiency of traditional dictionary-based compression techniques.
- (2) We show that multi-dictionary techniques (such as Ishiura and Yamaguchi [1997]) may have good theoretical compression ratio, but when actually implemented on FPGAs, they may occupy more memories than uncompressed code. To address this issue, we propose a technique that leverages on-chip dual-port memory blocks on FPGAs to improve memory utilization.
- (3) We also propose a multi-level (cascaded) dictionary architecture that can reduce the number of utilized on-chip memories by an additional 20% compared to single-level dictionary compression.

5. MULTI-DICTIONARY MICROCODE COMPRESSION

In this section, we describe the concept of multi-dictionary compression and explain how “don’t care” values in microcode can be leveraged for code size reduction. Our tool constructs a dictionary of unique control words and, in the executable binary, replaces each control word by its corresponding dictionary line addresses. Figure 5 shows a one-dictionary (*opt1*) code compression approach. The memory structure consists of a *Code lookup table (CodeLUT)* and a dictionary. The Program Counter (PC) contains the address of *CodeLUT* and is used to read the next *codeword*. The codeword is then used to read the corresponding control word from dictionary.

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	0	0	1	0	1	1	0	0	0	1	0	1	0	1	0
1	0	0	0	1	1	1	1	0	0	1	0	1	0	1	0
1	0	0	1	0	1	1	0	1	1	1	0	1	0	1	0
1	0	0	1	0	1	1	0	0	0	1	0	1	0	1	0
0	0	1	0	1	0	1	0	1	0	0	0	1	1	1	1
0	0	1	0	1	0	1	0	1	0	0	1	0	1	1	0
1	0	0	0	1	1	1	1	0	0	1	0	1	0	1	0
0	0	1	0	1	0	1	0	1	0	0	0	1	1	1	1
1	0	0	1	0	1	1	0	1	1	1	0	1	0	1	0

Fig. 6. CWs of a sample program.

000																
001																
010																
000	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
011	1	0	0	1	0	1	1	0	0	0	1	0	1	0	1	0
100	1	0	0	0	1	1	1	1	0	0	0	1	0	1	0	1
001	1	0	0	1	0	1	1	0	1	1	1	0	1	0	1	0
011	0	0	1	0	1	0	1	0	1	0	0	0	1	1	1	1
101	0	0	1	0	1	0	1	0	1	0	0	1	0	1	1	0
010																

Fig. 7. Single-dictionary compression on CWs of Figure 6.

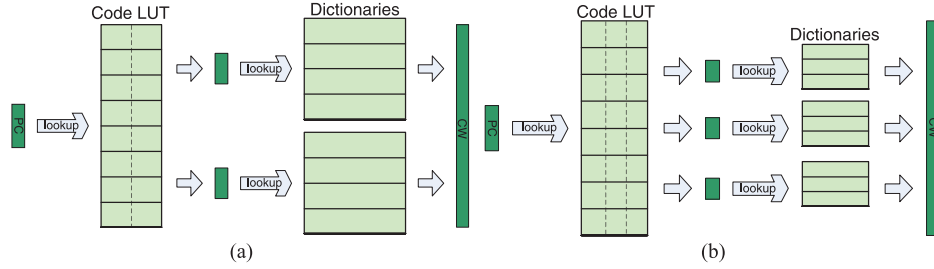
The following example shows how dictionary-based compression can reduce the total size of control memory. Suppose that Figure 6 shows the CWs of a sample program. Each CW has 16 bits and the program has nine CWs. Therefore, the code size of the program is 144 bits (16×9). Figure 7 shows the compressed implementation of Figure 6, where the dictionary contains five unique CWs and the CodeLUT contains the corresponding address of the CWs. To address the dictionary, three bits are needed; thus, the codewords are three-bit wide. After compression the total binary size is reduced to 107 (i.e. $3 \times 9 + 16 \times 5$).

Since CWs can be very wide with many unique patterns, the dictionary may have many entries and the compression efficiency may be low. To increase the chances of finding matching patterns, we can partition the CWs to smaller slices and construct multiple dictionaries. Usually the total size of the partitioned dictionaries is much smaller than that of a single big dictionary. However, corresponding to each dictionary, a code field must be added to codewords. Figure 8 shows the two-dictionary (*opt2*) implementation of the example CWs of Figure 6. The top dictionary contains the unique patterns of the least significant half of the CWs, while the bottom dictionary has those of the most significant halves. Note that the number of unique control word slices in each half is less than the total number of unique control words. The codewords have two fields, which are used to address the two dictionaries. Since each dictionary has four or less entries, the codeword fields are only two bits. Using two-dictionary implementation the code size is reduced to 92 bits (i.e. $4 \times 9 + 8 \times 3 + 8 \times 4$).

As number of dictionaries increases, the number of codeword fields increases and eventually cancels out the code size reduction achieved by partitioning.

00	00	0	1	2	3	4	5	6	7
01	00	1	0	0	1	0	1	1	0
00	01	1	0	0	0	1	1	1	1
00	00	0	0	1	0	1	0	1	0
10	10	8	9	10	11	12	13	14	15
10	11	0	0	1	0	1	0	1	0
01	00	1	1	1	0	1	0	1	0
10	10	1	0	0	0	1	1	1	1
00	01	1	0	0	1	0	1	1	0

Fig. 8. Two-dictionary compression on CWs of Figure 6.

Fig. 9. (a) Two-dictionary (*opt2*) and (b) Three-dictionary compression (*opt3*).

1	0	X	X	1	1	X	X	X	1	1	0	X	0	1
0	X	X	1	1	0	0	0	1	X	0	X	X	1	1
X	0	0	1	X	1	0	1	0	X	X	0	0	0	X
1	X	0	X	1	X	0	X	0	1	1	X	0	X	1
0	0	X	1	1	0	X	0	1	0	0	0	0	1	1

Fig. 10. Example of control words generated by NISC compiler.

Figure 9 shows two-dictionary (*opt2*) and three-dictionary (*opt3*) code compression approaches. The performance penalty in all cases is the same since the dictionaries are accessed in parallel. In addition to the number of dictionaries, the way “X” values are resolved in the binary may affect the efficiency of compression.

5.1 Compression-Aware “Don’t Care” Resolution (CX)

Figure 10 shows an example of NISC control words. As explained in Section 2, the control words contain “don’t care” values (denoted by “X”), which indicate that some of the units are idle at a given cycle and their control signals can be assigned to “0” or “1” without affecting program behavior. To build a dictionary for the CWs, one may replace “X” values by “0” and then extract the unique patterns. In that case, the dictionary (shown in Figure 11) will have four entries, because only the second and the last vectors match. However, if the “X” values are smartly resolved, then seemingly different patterns can be combined into one dictionary entry.

1	0	0	0	1	1	0	0	0	1	1	0	0	1
0	0	0	1	1	0	0	0	1	0	0	0	0	1
0	0	0	1	0	1	0	1	0	0	0	0	0	0
1	0	0	0	1	0	0	0	0	1	1	0	0	1

Fig. 11. Dictionary content for CWs of Figure 10 (“X” are replaced by “0”).

In general, we need to resolve “X” values in CWs so that the total number of unique patterns is minimized. To solve this problem, we convert it to a graph coloring [Jensen and Toft 1995] problem. For a given list of bit-vectors, we construct a graph $G(V, E)$, where the vertices in V are the bit-vectors, and the edges in E show the conflict between the vectors. Two bit-vectors v_1 and v_2 do NOT have conflict if they can be collapsed into a single vector:

$$\forall i \in \{1, \dots, N\}, v_1[i] = v_2[i] \text{ OR } v_1[i] = \text{“X”} \text{ OR } v_2[i] = \text{“X”}.$$

where, N is the number of bits in a bit-vector. The edges in E are defined between the vectors that have conflict with each other:

$$E = \{(v_1, v_2) \mid v_1 \text{ has conflict with } v_2\}.$$

The algorithm must partition the vertices (or vectors) to subcategories so that there is no edge (i.e., conflict) between any two vertices in the same category while minimizing the total number of categories. This is exactly the graph coloring problem where each category is represented by a distinct color [Jensen and Toft 1995]. Solving the graph coloring problem optimally is NP-hard [Garey and Johnson 1979]. But there are many well-known heuristics that generate fine results in polynomial time. We use Welsh and Powell algorithm [Jensen and Toft 1995], a greedy heuristic, that (1) sorts vertices based on their degree in decreasing order; (2) traverses the graph and colors as many nodes as possible with color c_1 ; and (3) repeats step (2) with color c_2 , c_3 , etc., until no vertex is left uncolored. Figure 12 shows the details of our algorithm. After coloring the graph, a new vector is generated for each color by combining the values of the vectors that share the same color (see line 11–16). Then, all such vectors are replaced with the new vector. The new vectors are also used to fill the dictionary. Figure 13 shows how the algorithm is applied to the example of Figure 10. The graph coloring algorithm produces only two colors. The first, third, and fourth vectors are mapped to the first entry of Figure 14, and the other two vectors are mapped to the second entry.

This algorithm is coded in C# and is added to our toolset. For our benchmarks, it takes only a few seconds to apply this algorithm. The only exception is MP3, which takes a few minutes.

In Gorjiara and Gajski [2007], we have shown that “X” values can also be resolved for power optimization at the cost of compromising compression efficiency. We also proposed a profile-driven “don’t care” resolution technique that achieves both power and code size efficiency.

```

1  Compress ( $V, N$ )
2  //Inputs:
3  //V: the list of vectors
4  //N: the bit width of the vectors
5  //output: CV, the list of compact vectors.
6   $G = \text{ConstructConflictGraph}(V)$ ;
7   $C = \text{Color}(G)$ ; // C is the set of colors
8  for each color  $c$  in  $C$ 
9    // create a new vector that contains only 'X'
10    $cv = \text{new BitVector}('X', N)$ ;
11   for each vector  $v$  in  $V$ 
12     if (color of  $v$  is  $c$ )
13       //merging the two vectors
14       for  $i=1$  to  $N$ 
15         if ( $v[i] \neq 'X'$ )
16            $cv[i] = v[i]$ ;
17 for each vector  $v$  in  $CV$ 
18   replace the remaining 'X' in  $v$  with '0'
19  $CV.\text{Add}(cv)$ ;
20 return  $CV$ ;

```

Fig. 12. Compression-aware “don’t care” resolution algorithm.

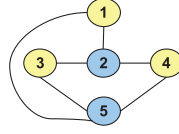


Fig. 13. Colored conflict graph of the CWs in Figure 10.

1	0	0	1	1	1	0	1	0	1	1	1	0	0	1
0	0	0	1	1	0	0	0	1	0	0	0	0	1	1

Fig. 14. Dictionary content using compression-aware “X” resolution.

5.2 Flow of Our Tool

Figure 15 shows the flow of our tool. First the application is compiled on the datapath and control words are generated. Then, the control words are partitioned and their “don’t care” values are resolved using the approach in Section 5.1. Finally, the HDL code of the design is generated. The techniques presented in this article have been completely implemented and are added to our toolset, which is available online at <http://www.cecs.uci.edu/~nisc>.

5.3 Compression Efficiency

In this section we study the effects of multi-dictionary compression as well as compression-aware “don’t care” resolution (CX). The same MP3 and MiBench benchmarks used in Section 3 are also used in this Section. Table III shows the binary size of the benchmarks compiled on GNISC with and without code compression. The second column (No-opt) shows the baseline code size without any compression. The third to sixth columns (opt1, opt2, opt3, and opt4) show

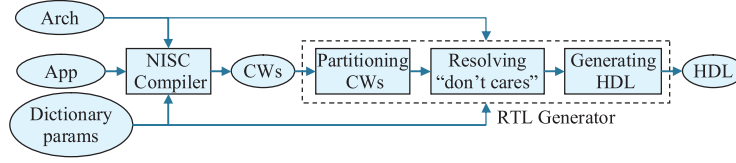


Fig. 15. Flow of our tools.

Table III. Code Size of Benchmarks with Different Compression Techniques

	without compression-aware 'x' resolution					with compression-aware 'x' resolution			
	No-opt	opt1	opt2	opt3	opt4	opt1	opt2	opt3	opt4
adpcm_coder	6.96	4.34	2.90	2.63	2.65	3.77	2.44	2.19	2.19
adpcm_decoder	5.08	3.76	2.38	1.92	2.03	3.24	1.98	1.59	1.68
CRC32	2.57	2.12	1.31	1.08	1.02	1.72	1.04	0.85	0.80
dijkstra	9.61	4.92	3.17	3.02	3.22	4.32	2.71	2.52	2.68
sha	14.12	8.24	5.61	5.23	5.26	6.75	4.45	4.12	4.14
Mp3	216.66	92.66	73.17	79.16	88.98	82.00	63.08	67.66	76.05
average CR	1.00	0.62	0.41	0.37	0.38	0.53	0.34	0.30	0.31

the code size with one to four dictionaries, respectively. In these approaches, the “X” values are simply replaced by “0” and then the dictionaries are constructed. As the number of dictionaries increases, the code size (i.e., the total size of dictionaries and *CodeLUT*) of all the benchmarks decreases up to certain points and then increases again. These are the points where the increase in *CodeLUT* size cancels out the benefit of having more dictionaries. The optimum number of dictionaries may vary for different applications.

The seventh to tenth columns of Table III show the binary size when dictionary-based compression is combined with the compression-aware “X” resolution (CX) technique introduced in Section 5.1. The CX technique reduces the code size by an additional 15%–27% (avg. 20%) compared to compression alone. Similarly, as the number of dictionaries increases, the code size of all the benchmarks decreases up to certain points (the highlighted values) and then increases again.

The last row of Table III shows the Compression Ratio (CR), a metric commonly used to evaluate a compression algorithm. CR is the ratio between the compressed size and the original size, and smaller CR values show a better compression. On average, for all these benchmarks, the CX-based three-dictionary compression (i.e., *opt3*) outperforms the others and achieves CR of 0.3. In *opt3*, the total code size is one-third of the code size of *No-opt*.

Table IV compares the code size of compressed GNISC with that of MicroBlaze. The second column shows the code size of optimized GNISC (i.e. the highlighted values in Table III) and the third column shows the code size of MicroBlaze (from Table II). The fourth column shows the ratio of the code sizes of GNISC and MicroBlaze. For the small benchmarks, the code size of compressed GNISC is very close or even smaller than that of MicroBlaze. However, for the medium and large benchmarks, GNISC code size is still significantly higher than MicroBlaze (30-40%). In Section 7, we propose a cascaded dictionary structure that can further reduce the code size of larger applications.

Table IV. Comparing Code Size of Compressed GNISC with MicroBlaze

	Code size (Kbytes)		Code size ratio
	GNISC-opt	MicroBlaze	GNISC-opt vs. MicroBlaze
adpcm_coder	2.19	1.95	1.12
adpcm_decoder	1.59	1.36	1.17
CRC32	0.80	1.26	0.63
dijkstra	2.52	1.93	1.31
sha	4.12	3.16	1.30
Mp3	63.08	44.62	1.41
average			1.16

Table V. Number of Utilized 18Kbit Block RAMs

	Memory size (number of 18Kbit block RAMs)					
	MBlaze	No-opt	opt1	opt2	opt3	opt4
adpcm_coder	1	4	5	6	6	7
adpcm_decoder	1	3	4	5	5	6
CRC32	1	3	4	5	5	6
dijkstra	1	6	5	6	7	9
sha	2	11	5	6	9	11
Mp3	21	117	67	34	38	38

In this section, we showed that multi-dictionary compression along with compression-aware “X” resolution reduce the code size by more than three times (i.e., compression ratio of 0.3). However, in the next section we show that low compression ratio does not necessarily result in an efficient design when targeting FPGA platforms. In fact, in some cases, the compressed code may occupy more resources than the uncompressed one.

5.4 Block RAM Utilization in FPGA-Based Implementations

In this section, we investigate whether dictionary-based compression can actually reduce resource utilization in an FPGA-based implementation. In FPGAs, the *CodeLUT* and dictionaries may be implemented using lookup tables or memory blocks (RAM). In today’s FPGAs, tens or even hundreds of compact and fast memory blocks exist. Each block has a predefined size and a set of configurations: for example, in Xilinx Virtex4 FPGA, each block RAM is 18Kbits and can be configured statically to a $1 \times 18\text{Kb}$, $2 \times 9\text{Kb}$, $4 \times 4\text{Kb}$, $8 \times 2\text{Kb}$, $16 \times 1\text{Kb}$, or $32 \times 0.5\text{Kb}$ memory. These configurations are called RAM *primitives*. In FPGAs, logical memories are implemented using one or more block RAMs depending on their width, depth, and available primitives. Reducing the number of utilized BRAMs is important because it allows packing more processing elements into smaller, low-cost FPGAs.

Table V shows the number of utilized 18Kbit block RAMs in different implementations on Xilinx Virtex4SX35. This package is large and contains hundreds of block RAMs. In the MicroBlaze implementation (the second column), most of the benchmarks need only one block RAM for their code, except for *sha* and *Mp3*, which need 2 and 21 blocks, respectively. These numbers are significantly higher for NISC (the third column), because the CWs are wide, and block RAM primitives do not support wide words. In terms of block RAM utilization, NISC requires on average five times more blocks than MicroBlaze.

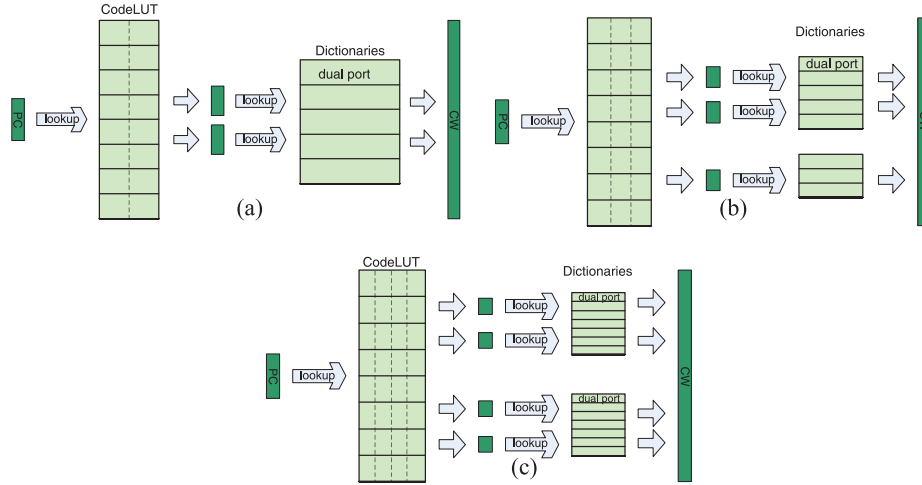


Fig. 16. (a) Two-, (b) three-, (c) four-dictionary code compression using a dual-port memory (*opt2DP*, *opt3DP*, *opt4DP*).

Surprisingly, the compression techniques even increase the number of utilized block RAMs for the smaller applications (i.e., *adpcm_coder*, *adpcm_decoder*, and *CRC32*). However, for medium and large applications (i.e., *dijkstra*, *sha*, and *Mp3*), the compression techniques reduce the number of block RAMs. Although our code compression techniques reduce the code size, they tend to increase the number of memory units, thus occupying more partially-utilized block RAMs. Note that *opt3* is the best compression technique in terms of compression ratio (as shown in Table III). However, it wastes many block RAMs in FPGA implementations (as shown in Table V). To address this issue, we introduce our FPGA-aware microcode compression in the following section.

6. FPGA-AWARE MICROCODE COMPRESSION

In Xilinx FPGAs, each block RAM can be configured as single port or dual port with very little logic overhead. We use this property to reduce the number of utilized block RAMs. We integrate every two dictionaries into one dual-port memory. Figure 16(a), (b), and (c) show the dual-port implementation of *opt2*, *opt3*, and *opt4*, respectively. Note that the merged dictionary contents may have more entries than each individual dictionary. However, the size of the merged dictionary is less than the total size of the two dictionaries, because redundant entries can be removed after merging the contents. Since merging dictionaries increases the depth of the dictionary unit, the width of codewords may increase as well. As a result, the total code size most-likely remains the same as before, but the number of utilized block RAMs decreases.

Figure 17 shows the dictionary content and *CodeLUT* of the dual-port implementation of Figure 8. The seven entries of the two dictionaries in Figure 8 are compacted to four unique entries in Figure 17. The codewords are also updated to refer to the correct bit patterns. Compared to Figure 8 that requires

11	00
10	00
11	01
11	00
00	10
00	11
10	00
00	10
11	01

0	0	1	0	1	0	1	0
1	1	1	0	1	0	1	0
1	0	0	0	1	1	1	1
1	0	0	1	0	1	1	0

Fig. 17. Dual-port memory implementation of Figure 8.

Table VI. Number of Utilized 18Kbit Block RAMs for All Compression Approaches

	Memory size (number of 18Kbit block RAMs)								
	MBlaZe	No-opt	opt1	opt2	opt3	opt4	opt2DP	opt3DP	opt4DP
adpcm_coder	1	4	5	6	6	7	3	4	4
adpcm_decoder	1	3	4	5	5	6	2	3	3
CRC32	1	3	4	5	5	6	2	3	3
dijkstra	1	6	5	6	7	9	3	5	5
sha	2	11	5	6	9	11	4	5	5
Mp3	21	117	67	34	38	38	36	42	43

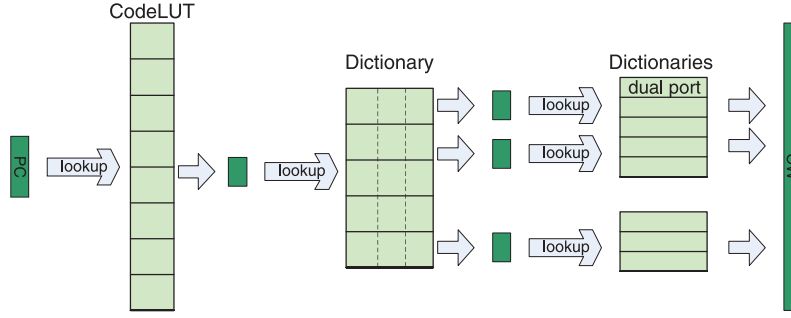
three block RAMs, Figure 17 requires only two RAMs. Also, the code size is reduced to 68 bits (i.e. $4 \times 9 + 8 \times 4$).

Table VI shows the number of utilized block RAMs with single and dual port compression techniques. Using dual-port memories (i.e., *opt2DP*, *opt3DP*, and *opt4DP*) reduces the number of utilized block RAMs compared to single-port memories. The minimum number of blocks is achieved using *opt2DP* for *adpcm_coder*, *adpcm_decoder*, *CRC32*, *dijkstra*, and *sha*. For *Mp3* decoder, however, the minimum is achieved using *opt2*. That is due to a significant size increase in *CodeLUT* of *MP3* when using dual-port dictionaries. Overall, merging the dictionaries (i.e., *opt2DP*) reduces the number of block RAMs by 46% compared to uncompressed NISC, and by 50% compared to corresponding unmerged dictionary compression (i.e. *opt2*). However, compared to MicroBlaze, *opt2DP* requires 2.3 times more Block RAMs, on average. Note that as the program size increases, the gap between *opt2DP* and MicroBlaze decreases.

To further improve RAM utilization, we can also adjust the width of dictionaries to match the width of RAM primitives. For example, in a three-dictionary implementation of 75-bit control words, instead of having three equally-sized 25-bit dictionaries, we can have two 32-bit wide dictionaries (that can be merged into one dual-port memory), and an 11-bit dictionary. This way, the block RAMs are more efficiently utilized and deeper memories can be designed with fewer RAMs.

7. CASCADING DICTIONARIES

For large applications, we also propose to use *cascaded dictionaries*, where the *CodeLUT* is replaced by a new dictionary and a narrower *CodeLUT*. Figure 18 shows the cascaded version of the *opt3DP* (see Figure 16(b)). In this approach, the unique three-field codewords are stored in an intermediate dictionary and *CodeLUT* content is replaced by the references to those dictionary

Fig. 18. Cascaded version of *opt3DP* in Figure 16(b).

entries. The new codewords must be narrower than the original codewords to achieve savings. This technique requires an additional lookup and hence has higher performance penalty. For small and medium size applications, cascading dictionaries improves the compression ratio but does not reduce the number of utilized Block RAMs. However, for the Mp3 application *cascaded-opt3DP* reduces the number of utilized block RAMs as well (from 34 to 28). If we also adjust the dictionary width to match Xilinx Block RAM primitives, we can reduce the number of blocks to 27, which is 4.3 times better than NISC without compression (No-opt), and 20% better than the best single-level compression technique for MP3 (i.e., *opt2*). Compared to MicroBlaze that needs 21 Block RAMs, cascaded, merged dictionary NISC needs only 28% more RAMs.

8. PERFORMANCE PENALTY OF DECOMPRESSION

As mentioned in Section 4, the performance penalty of a decompression unit depends on its complexity as well as its position with respect to the cache architecture. If the decompression unit is placed between the main memory and the cache (a pre-cache approach), then decompressions are limited to cache misses only and hence, their overall penalty is negligible. However, if the unit is placed between the cache and the processor (a post-cache approach) or between the memory and the processor (in systems without a cache) then, the penalty is non-negligible. In this section, to determine the worst-case performance penalty, we consider a system without a cache where decompression unit is placed between memory and the processor.

All the compression techniques in Sections 5 and 6 have the same performance penalty. They increase the number of fetch pipeline stages by one, which increases the branch delay by one cycle. The cascaded dictionary structure proposed in Section 7 has a higher performance penalty because it increases fetch pipeline stages (and hence the branch delay) by two. In addition to caches, there are two other approaches for reducing the performance penalty of decompression: (1) branch prediction, and (2) filling branch delay slot. In branch prediction, the controller predicts what would be the target of a branch operation, and starts fetching from the predicted target address. If it turns out that the prediction was wrong, the controller flushes the affected pipeline stages. In

Table VII. Comparing Performance of GNISC-opt with That of GNISC and MicroBlaze

	GNISC (without compression)	GNISC-opt. (with 1-level compression)	GNISC-opt. vs. GNISC	GNISC-opt. vs. MicroBlaze
Benchmarks	#cycles	#cycles	slowdown (%)	speedup (x)
adpcm_coder	74321930	84251684	13.36	3.05
adpcm_decoder	63082673	66504319	5.42	4.85
CRC32	21901993	26008604	18.75	8.05
dijkstra	9764682	10631310	8.88	2.44
sha	19282976	18371827	3.33	9.96
Mp3	897452	927307	4.96	2.88
Average			9.12	5.21

the branch delay slot approach, the compiler tries to find operations that are independent of the branch operation, and schedule them after the branch to fill the branch delay slot. The branch delay slot approach is done at compile time and does not impose any hardware overhead. In contrast, branch prediction needs more complex hardware. Currently, our compiler supports branch delay slot, but branch prediction can also be added in the future.

We generated HDL code for the GNISC datapath with and without compression stage(s) and simulated the code using Modelsim simulator. In Table VII, the second and third columns show cycle count of each benchmark without and with compression, respectively. The fourth column shows the performance overhead of single-level compression in terms of the slowdown percentage compared to the baseline GNISC. The performance penalty depends on the number of jump operations and how well the compiler can fill the extra branch delay slot. On average, the performance is degraded by only 9.12% (up to 19%). The fifth column compares the speed of optimized GNISC with that of MicroBlaze processor. The optimized GNISC is on average 5.21 times faster than MicroBlaze. This shows that the compression techniques had little effect on the performance of GNISC even without cache and branch prediction. For cascaded dictionaries (two-level compression), only the MP3 application shows improvement. Therefore, we simulated MP3 and observed an additional 3% performance penalty compared to single-level compression. This penalty is about 8% when compared to uncompressed GNISC.

9. CONCLUSION

In this article, we studied the code size of NISC PEs and compared it with that of traditional RISC processors. We observed that although NISC PEs outperform RISC processors by five times on average, their code sizes are about four times larger.

We studied the use of different variations of dictionary-based code compression techniques on the NISC binary. We showed that although multi-dictionary compression reduces the code size by 3.3 times, its FPGA implementation is very inefficient and can result in occupying more resources than the uncompressed code. To overcome this limitation, we proposed to merge every two dictionaries into a single dual-port memory unit on FPGAs. Using this approach, the block RAM utilization is improved by 46%. Also, for large applications, we proposed using *cascaded dictionaries*, where multi-levels of dictionaries are used to decompress the code. For MP3 application, a merged, cascaded,

three-dictionary implementation reduces the number of utilized block RAMs by 4.3 times (76%) compared to a NISC without compression. This corresponds to 20% additional savings over the best single-level dictionary-based compression. Compared to MicroBlaze our compressed MP3 consumes only 28% more block RAMs.

REFERENCES

- AGRAWALA, A. AND RAUSCHER, T. 1976. *Foundations of Microprogramming: Architecture, Software, and Applications*. Academic Press.
- CODWELL, R., NIX, R., DONNELL, J., PAPWORTH, D., AND RODMAN, P. 1987. A VLIW architecture for a trace scheduling compiler. *ACM SIGOPS Operat. Syst. Rev.* 21, 4.
- CORLISS, M., LEWIS, E., AND ROTH, A. 2003. DISE: a programmable macro engine for customizing applications. In *Proceedings of the International Symposium on Computer Architecture (ISCA)*.
- FRASER, C. 2002. An instruction for direct interpretation of LZ77-compressed programs. Tech. rep. MSR-TR-2002-90, Microsoft Research, Microsoft Corporation.
- GAREY, M. AND JOHNSON, D. 1979. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W.H. Freeman.
- GORJIARA, B. 2007. Synthesis and optimization of custom low-power NISC processors. Ph.D. dissertation, University of California, Irvine.
- GORJIARA, B. AND GAJSKI, D. 2005. Custom processor design using NISC: A case-study on DCT algorithm. In *Proceedings of the IEEE Workshop on Embedded Systems for Real-Time Multimedia (ESTIMedia)*.
- GORJIARA, B. AND GAJSKI, D. 2007. A novel profile-driven technique for simultaneous power and code-size optimization of nanocoded IPs. In *Proceedings of the International Conference on Computer Design (ICCD)*.
- GORJIARA, B. AND GAJSKI, D. 2008. Automatic Architecture Refinement Techniques for Customizing Processing Elements. In *Proceedings of the Design Automation Conference (DAC)*.
- GORJIARA, B., RESHADI, M., CHANDRAIAH, P., AND GAJSKI, D. 2006. Generic netlist representation for system and PE level design exploration. In *Proceedings of the International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*.
- GREHAN, R. 1997. 16-bit: The good, the bad, your options. *Embed. Syst. Prog.*
- ISHIURA, N. AND YAMAGUCHI, M. 1997. Instruction code compression for application specific VLIW processors based on automatic field partitioning. In *Proceedings of the International Conference on Synthesis and System Integration of Mixed Information System (SASIMI)*.
- JENSEN, T. AND TOFT, B. 1995. *Graph Coloring Problems*. Wiley-Interscience. New York.
- KEMP, T., MONTOMEY, R., AUERBACK, D., HARPER, J., AND PALMER, J. 1998. *A Decompression Core for PowerPC*. IBM Corporation.
- LAU, J., SCHOENMACKERS, S., SHERWOOD, T., AND CALDER, B. 2003. Reducing code size with echo instructions. In *Proceedings of the International Conference on Compilers, Architecture and Synthesis for Embedded Systems (CASES)*.
- LEFURGY, C., PICCININNI, E., AND MUDGE, T. 1999. Evaluation of a high performance code compression method. In *Proceedings of the International Symposium on Microarchitecture*.
- LEKATSAS, H., HENKEL, J., AND JAKKULA, V. 2002. Design of a one-cycle decompression hardware for performance increase in embedded systems. In *Proceedings of the Design Automation Conference (DAC)*.
- PRAKASH, J., SANDEEP, C., SHANKAR, P., AND SRIKANT, Y. 2003. A simple and fast scheme for code compression for VLIW processors. In *Proceedings of the Data Compression Conference*.
- RAFAIL, K. 1994. *Universal Compression and Retrieval*. Kluwer Academic. Publishing.
- RAU, B., YEN, D., YEN, W., AND TOWLE, R. 1989. The cydra 5 departmental supercomputer: Design philosophies, decisions, and trade-offs. *IEEE Computers*, 22, 1, 12–34.
- ACM Transactions on Reconfigurable Technology and Systems, Vol. 1, No. 2, Article 11, Pub. date: June 2008.

- RESHADI, M. 2007. No-instruction-set-computer (NISC) technology modeling and compilation. Ph.D. thesis, University of California, Irvine.
- RESHADI, M. AND GAJSKI, D. 2005. A cycle-accurate compilation algorithm for custom pipelined datapaths. In *Proceedings of the International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*.
- RESHADI, M. AND GAJSKI, D. 2007. Interrupt and low-level programming support for expanding the application domain of statically-scheduled horizontally-microcoded architectures in embedded systems. In *Proceedings of the Design Automation and Test in Europe (DATE)*.
- RESHADI, M., GORJIARA, B., AND GAJSKI, D. 2005. Utilizing horizontal and vertical parallelism using a no-instruction-set compiler and custom datapaths. In *Proceedings of the International Conference on Computer Design (ICCD)*.
- RESHADI, M., GORJIARA, B., AND GAJSKI, D. 2008. C-Based design flow: A case study on G.729A for voice over internet protocol (VoIP). In *Proceedings of the Design Automation Conference (DAC)*.
- ROS, M. AND SUTTON, P. 2004. A hamming distance based VLIW/EPIC code compression technique. In *Proceedings of the International Conference on Compilers, Architectures, and Synthesis for Embedded Systems (CASES)*.
- SAGHIR, M. 1998. Application-specific instruction-set architectures for embedded SDP applications. Ph.D. thesis, University of Toronto.
- SEGARS, S., CLARKE, K., AND GOUDGE, L. 1995. Embedded control problems, Thumb, and the ARM7TDMI. *IEEE Micro* 15, 5, 22–30.
- TRAJKOVIC, J., RESHADI, M., GORJIARA, B., AND GAJSKI, D. 2006. A graph based algorithm for data path optimization in custom processors. In *Proceedings of the Euromicro Conference on Digital System Design*.
- WANG, K. 2001. Code compaction for VLIW instructions. M.S. thesis, University of Toronto.
- WEBER, S. AND KEUTZER, K. 2005. Using minimal minterms to represent programmability. In *Proceedings of the International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*.
- WOLFE, A. AND CHANIN, A. 1992. Executing compressed programs on an embedded RISC architecture. In *Proceedings of the International Symposium on Microarchitecture*.
- XIE, Y., WOLF, W., AND LEKATSAS, H. 2001. A code decompression architecture for VLIW processors. In *Proceedings of the International Symposium on Microarchitecture*.
- XIE, Y., WOLF, W., AND LEKATSAS, H. 2001. Compression ratio and decompression overhead tradeoffs in code compression for VLIW architectures. In *Proceedings of the IEEE International ASIC Conference*.
- XIE, Y., WOLF, W., AND LEKATSAS, H. 2002. Code compression for VLIW processors using variable-to-fixed coding. In *Proceedings of the International Symposium on System Synthesis (ISSS)*.
- XIE, Y., WOLF, W., AND LEKATSAS, H. 2003. Profile-driven selective code compression. In *Proceedings of the Design, Automation and Test in Europe (DATE)*.

Received May 2007; revised October 2007, February 2008; accepted April 2008