Compiling For SIMD Within A Register

Randall J. Fisher and Henry G. Dietz

School of Electrical and Computer Engineering
Purdue University
West Lafayette, IN 47907-1285
rfisher@ecn.purdue.edu and hankd@ecn.purdue.edu

Abstract

Although SIMD (Single Instruction stream Multiple Data stream) parallel computers have existed for decades, it is only in the past few years that a new version of SIMD has evolved: SIMD Within A Register (SWAR). Unlike other styles of SIMD hardware, SWAR models are tuned to be integrated within conventional microprocessors, using their existing memory reference and instruction handling mechanisms, with the primary goal of improving the speed of specific multimedia operations.

Because the SWAR implementations for various microprocessors vary widely and each is missing instructions for some SWAR operations that are needed to support a more general, portable, high-level SIMD execution model, this paper focuses on how these missing operations can be implemented using either the existing SWAR hardware or even conventional 32-bit integer instructions. In addition, SWAR offers a few new challenges for compiler optimization; these issues are briefly introduced, as are the SWARC module language and the Scc compiler.

1. Introduction

The SWAR model takes advantage of the fact that a wide data path within a processor also can be treated as multiple, thinner, SIMD-parallel data paths. Thus, each register is effectively partitioned into *fields* that can be operated on in parallel. Operations like partitioned field-by-field addition are easily implemented with minimal hardware enhancements. For example, to make a 64-bit adder function as eight 8-bit adders, one simply needs to modify the adder's carry logic so that carry from one 8-bit field to the next is suppressed. But what SWAR field sizes and operations should be supported?

Because the introduction of SWAR techniques was motivated by the need to improve multimedia application performance, it is not surprising that the field sizes supported in hardware tend to be those commonly associated with multimedia's "natural" data types. In particular, 8-bit pixel color channel values are used for most forms of video and image processing, so this data size becomes very important. Similarly, 16-bit audio samples are very commonly used, and this field size also can be used for intermediate results when computations on 8-bit values need more than 8-bit precision. Other data sizes are common in other applications yet not directly supported by any of the SWAR hardware extensions in current microprocessors. For example: 1-bit fields hold boolean values, 2-bit fields hold either 0/1/X

values for a logic simulator or a base-pair value for operating on a gene database, and 10-bit or 12-bit values are used for a wide range of digitizing devices.

A high-level SWAR language can and should allow users to specify the minimum number of bits needed for each datum, without requiring that precise field size or implying a specific internal layout. As discussed in section 4, the SWARC module language and Scc compiler make use of these semantics by selecting the most efficient SWAR field size and layout that is available on a given target machine. Thus, different targets may use different layouts without violating the semantics assumed by the high-level language SWAR application programmer.

The operations implemented in SWAR hardware also are strongly biased in favor of specific multimedia applications and algorithms. Addition, subtraction, and various other integer operations are commonly supported. Because adding two bright pixel values should not result in a dark pixel, there are a variety of *saturation arithmetic* operations supported (e.g., overflow in saturation arithmetic yields the maximum representable value rather than the low bits of a larger value). 32-bit floating-point operations are also supported by AMD's 3DNow! and Motorola's AltiVec extensions.

Unfortunately, at least from the point of view of building high-level language compilers, current SWAR hardware support is limited to just a few field sizes and operations. Even worse, the supported operations and partitioning reflect the designers' expectations of the needs of multimedia programmers and the internal structures within the particular processor architectures, and are therefore inconsistent across architecture families.

While multimedia extensions were intended for specific applications, a generalized SWAR model may be used to create large code sequences which are not easily written, ported, or tuned by hand. A primary goal of this research, and of the SWARC module language, is to provide a consistent programming model so that writers of general-purpose SIMD programs can more easily take advantage of current and future SWAR capabilities. The following table illustrates the variance in multimedia support by showing the partitioning supported by each of five enhanced architecture families for typical arithmetic operations.

In table 1, the columns represent the multimedia extension set families. The extension families included are Digital Equipment Corporation extensions for the Alpha [DEC96]; Hewlett-Packard PA-RISC MAX v2.0 [Lee95][LeH96]; Silicon Graphics MIPS MDMX [SGI96]; Intel [Int96], Advanced Micro Devices [AMD97], Cyrix [Cyr96], and IDT WinChip [CeT97] versions of MMX; Sun SPARC V9 VIS [Sun97]; Motorola [Mot98] AltiVec; Cyrix [Cyr98] Extended MMX; and AMD [AMD98] 3DNow!.

The first few rows of the table describe the basic characteristics of the multimedia (MM) registers used for SWAR: How many registers are there? How wide is each register? Are these registers the processor's integer registers or are they overlaid on the floating-point register file? Are memory operands allowed or is a strict load/store instruction set used? Is floating point SWAR supported? The remaining rows sample various parallel arithmetic operations which one or more of the families support. "Part" is used to identify operations using partitioned register

operands, with the notation *axb*u specifying that the register holds *a* fields, each containing a *b*-bit value, with the u suffix indicating that the field values are unsigned and an f suffix indicating floating-point field values. "Single" refers to a single, scalar, operand value; "Immed" refers to an immediate value.

The many differences between these SWAR extension implementations are significantly more complex to manage than the variations in the base instruction sets of these processors. In order to execute fairly general SIMD code, there are also a number of instructions that are conspicuously absent from most SWAR instruction sets — for example, instructions to help implement SIMD-style conditional branches based on ANY or ALL condition tests. To support a high-level language, these problems must be efficiently managed by the compiler.

One obvious shortcoming of the SWAR hardware support is that it generally provides operations for just a few field sizes. Fortunately, we are not limited to the set of partitionings directly supported by the hardware. Logically, partitioning is not a property of "partitioned registers," but of an operation applied to the data in one or more registers. In other words, the data path through the function units may be partitioned in various ways, but the registers really are not partitioned; for example, there is nothing to prevent the programmer from treating a register's contents as having 8-bit fields in one operation and 16-bit fields in the next. It is even possible, and sometimes desirable, for operations to treat a register as though it is partitioned into field sizes which are not supported directly by the enhanced hardware. In fact, the partitioning need not even keep all fields within a register the same size, although managing irregular field sizes is too complex to be discussed in detail in this paper.

In summary, the fact that multimedia registers do not store partitioning information offers great flexibility. In a traditional SIMD execution model, data storage is allocated to each processing element, thus fixing both the number of data that can be operated on in parallel and the boundaries between values in different processors. Using SWAR, we can explicitly trade precision for parallelism width, using only as many bits for each operation as is algorithmically required. Further, the penalty for data crossing these imaginary inter-processing-element boundaries is zero for a SWAR system, which gives SWAR interesting new abilities for operations like reductions.

Our approach to making use of SWAR is based on creating a SWAR module language and compilers, as discussed in section 4, so that users can write full-featured portable SIMD programs that can be compiled into efficient SWAR modules and interface code that allows these modules to be used within ordinary C programs. It is not particularly difficult to build an optimizing SIMD module compiler supporting only those operations directly provided by a particular SWAR hardware implementation, but portability across all flavors of SWAR-capable hardware makes basic coding of constructs an interesting problem. Section 2 discusses the basic coding techniques needed to achieve efficient execution of general SIMD constructs on any type of SWAR hardware support — including 32-bit integer instruction sets that have no explicit support for SWAR. Beyond the basic coding of SWAR constructs, the SWAR execution model makes it appropriate to consider several new types of compiler optimizations, which are

	DEC Alpha	HP PA-RISC MAX	SGI MIPS MDMX	Intel, etc. MMX	Sun Sparc V9 VIS	Motorola AltiVec	Cyrix XMMX	AMD 3DNow!
# MM Registers	32	31	32/1	8+mem	32	32	8+mem	8+mem
# Bits/ Register	64	32/64	64/192	64	64	128	64	64
Which Regs?	Integer	Integer	Float	Float	Float	Vector	Float	Float
Mem->Reg?				•			•	•
Float Support?						•		•
Modular Add Part/Part 16		•	•	•	•	•	•	•
Saturation Add Part/Part 8u Part/Part 16 Part/Part 32		•	•	•		•	•	•
Sum of Diffs	•				•		•	
Modular Mul Part/Part 8u Part/Part 16 Part/Part 32f			•	•	•	•	•	•
Saturation Mul Part/Part 16			•					
1/x 32f								•
Shifts Part/Part 8u Part/Part 16u Part/Part 32u Part/Part 64u	•	•	•	•		•	•	•
Maximum Part/Part 8u Part/Part 32f	•		•			•		•
AND Part/Part	•		•	•		•	•	•
Compares Cond. Branch Cond. Load	=,<,≥,≤		=,<,≤	=,>		=,>,ALL	=,>	=,>,≥

Table 1: Partial Comparison Of Multimedia Instruction Set Extensions

briefly discussed in section 3. Section 4 unites all these concepts, describing the SWARC module language and the Scc compiler. Conclusions are summarized in section 5.

2. Basic SIMD-to-SWAR Compilation

In determining how to code high-level SIMD constructs using SWAR instructions, it is useful to distinguish several different classes of SIMD operations based on how they can be implemented using SWAR:

A polymorphic operation is a SIMD operation for which the same instruction(s) can be used to
implement the operation, independent of the field partitioning. For example, bitwise
operations like NOT, AND, OR, and XOR are polymorphic, as is the classical SIMD conditional

test ANY (which returns "true" if any bit in any field is a 1).

- A partitioned operation is a SIMD operation in which the value computed by each processing element is a function only of data that also resides in that processing element. In a SWAR implementation, a partitioned operation requires corresponding operand fields to be manipulated without interfering with adjacent fields. This is best accomplished using hardware partitioning, and SWAR hardware extensions generally support various partitioned operations, such as field addition and subtraction. However, no current SWAR hardware implements fully general partitioning for all operations, so software techniques to construct appropriate partitioned operations are critical.
- A communication operation logically transmits data across processing elements in an arbitrary pattern. For example, the MasPar MPL construct router[a].b accesses the value of the variable b on processing element a. Unfortunately, except for HP PA-RISC MAX and Motorola AltiVec, current SWAR extensions do not allow such general communication patterns to be used in re-arranging field values, and even MAX only allows this rearrangement within a single register (thus not solving the problem when SIMD vector lengths exceed one register). In fact, most SWAR hardware does not explicitly provide any communication operations; instead, they must be constructed using operations that can cross field boundaries, such as unpartitioned shifts.
- A type conversion operation converts SIMD data of one type into another. This becomes
 complicated in SWAR systems because data types that have different sizes yield different
 parallelism widths (numbers of fields per register). However, current SWAR hardware
 supports a variety of type conversion operations, and others can be constructed easily using
 communication and partitioned operations.
- A reduction operation recursively combines field values into a single value. In a traditional SIMD architecture, this involves communication between a shrinking set of processing elements working on values slowly growing in precision; this is a much more natural procedure for SWAR, and is thus worth exploring.
- A *masking* operation allows some processing elements to disable themselves. Obviously, for SWAR, there is no way to disable part of a register; however, there are arithmetic techniques that can be used to achieve the same result.
- A control flow operation is essentially a branch instruction that is executed by the SIMD control unit. Although SWAR systems do not have a control unit per se, the ordinary processor instruction set can be viewed as providing this functionality. In order to simplify pipeline structure, most SWAR hardware does not directly allow branching based on the contents of a partitioned SWAR register; it is thus necessary to move the SWAR fields into an ordinary register and test it there. These sequences are often awkward, but are unsurprising.

The following sections describe the basic coding of each of the more interesting classes: partitioned, communication, reduction, and masking operations.

2.1. Partitioned Operations

Partitioned operations are the primary focus of most of the hardware support for SWAR. In particular, most of the speedup claims for SWAR are based on partitioned additions of 8-bit fields, etc. However, not all important field sizes are handled by all SWAR hardware, and many partitioned instructions are simply omitted. We describe three different methods for implementing partitioned operations.

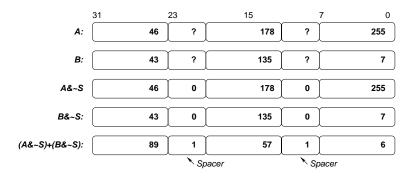
2.1.1. Hardware-Partitioned Operations

In an ideal partitioned operation, the operation is applied both concurrently and independently to the entire set of register fields. Thus, where SWAR hardware supports this, the SWAR code looks just like a single pure SIMD instruction. For example, consider adding two 4x8 values, A and B (each stored in a 32-bit register), using SWAR partitioned unsigned modular addition:

	31	23	15	7 0
A:	128	46	178	255
B:	49	43	135	7
A+B:	177	89	57	6

Numbering the fields from right to left, the addition performed in field 3 is 128+49 and yields the field result 177. Because this value is storable in an 8-bit field, it is not modified for storage. In contrast, the addition performed in field 1 is 178+135 and yields 313. This value requires nine bits for proper storage, thus it is modularized to fit in the eight bits which are available. The result stored for field 1 is then 313%(2⁸), or 57. Note that the overflow bit is lost and does not interfere with the operation as performed in field 2.

2.1.2. Partitioned Operations Using Spacer Bits



When an ideal partitioned addition is not supported by the hardware, spacer bits may be used to allow the field operations to be applied concurrently and independently (as shown above). These spacer bits are used as buffer zones between data fields, so that overflow and/or borrow can occur without interfering with adjacent fields. The spacer bit technique may be used not only to implement additional field sizes for SWAR hardware, but also to implement SWAR

partitioned operations on architectures providing only conventional 32-bit integer operations.

The example shown is essentially the same one used in section 2.1.1, but use of spacer bits allows only 3x8 values to fit in each register. At the start of the operation, the spacer values are unknown. This is indicated by a question mark (?) in each spacer field. To ensure that none of the field additions will overflow into the next field, the spacer bits of the addends are preset, or normalized, to zero before the addition is performed. This is done by ANDing the addends with a mask S, which has 1 bits only in the spacer positions.

The addition performed in field 2 occurs just as for an ideal partitioned addition, with no overflow and no modification of the stored result. The addition performed in field 1 is 178+135, yielding 313. This value requires nine bits for storage, and is stored with its lower eight bits in field 1, and its ninth bit carried into the spacer between fields 1 and 2. Only the part of the result stored in field 1 is considered to be valid. Thus, the valid result stored is $313\%(2^8)$ or 57. Note that the overflow bit does not interfere with the operation as performed in field 2. Similarly, the addition in field 0 results in a carry to the spacer between fields 0 and 1, and the storage of the result 6. Spacer bits also can be used to implement partitioned subtraction using unpartitioned instructions.

Only one spacer bit is needed between fields for addition or subtraction, but use of multiple spacer bits between the fields may allow multiple partitioned operations to be performed before re-normalizing the spacer bits. This static optimization simply requires tracking the range of possible values of the spacer bits to determine when re-normalization would be required.

2.1.3. Software-Partitioned Operations

Although the spacer-based partitioned operation code is fast, it only allows 3x8 values in each 32-bit register — not the 4x8 that could be used with hardware partitioning. With a few more instructions, the full densely-packed partitioned operation can be implemented.

The trick is simply to consider each field as being one bit smaller than it truly is, replacing the most significant bit of each field by a "virtual" spacer bit. After performing the spacer-partitioned operation on the modified register value, the most significant bits of each field are computed and inserted in the result. Thus, adding two 4x8 values is done as shown on the following page.

At the start of the operation, the two addends are each split into two partitioned registers. $\tt C$ and $\tt D$ contain the data — with the virtual spacer bit positions cleared in preparation for the partitioned addition. The resulting value for $\tt E$ is the correct 4x8 value, except in that the most significant bits of fields in $\tt A$ and $\tt B$ have not been added in. Fortunately, the bitwise $\tt XOR$ operation implements a one-bit add with no carry; thus, we can compute the addition of the most significant bits of each field by $\tt XOR$ ing the most significant bits of $\tt A$, $\tt B$, and $\tt E$ (with appropriate masking of the other bits).

	31	23		15		7	0
A:	12	8	46		178		255
B:	4	9	43		135		7
C=A&~S:	0	0 0	46	0	50	0	127
D=B&~S:	0 4	9 0	43	0	7	0	7
E=C+D:	0 4	9 0	89	0	57	1	6
F=A^B:	1 4	9 0	5	0	53	1	120
G=F&S:	1	0 0	0	0	0	1	0
H=E^G:	1 4	9 0	89	0	57	0	6
Н:	17	7	89		57		6

Note that the result field values in H are the same as those for ideal partitioned addition.

2.2. Inter-Processing-Element Communication

	31	23	15	7 0
A0:	3	2	1	0
A1:	7	6	5	4
B=A0>>8:	0	3	2	1
C=A1<<24:	4	0	0	0
A0'=B C:	4	3	2	1
A1'=A1>>8:	0	7	6	5

Partitioned registers may allow mesh and toroidal communication to be performed by applying shift and rotate instructions to the register. In the example above, A1 and A0 are, respectively, the upper and lower halves of an 8x8 vector layed out across two 4x8 partitioned registers. Although the numbering of fields is somewhat arbitrary, consider a right-neighbor communication as each processing element sending its value one field to the right.

Notice that the new value A0' was harder to compute than A1' because the datum from processing element 0 logically fell off the right side of the machine. If the right-neighbor communication had included a wrap-around (toroidal) link, the computation of A1' would have been much more similar to that of A0'.

Unfortunately, more complex communication patterns are not directly supported by most of the SWAR implementations. Implementations of more general communication patterns can be built using the PACK and UNPACK operations of SWAR extensions like MMX to implement shuffle and inverse-shuffle communication patterns that can be composed to simulate a

multistage interconnection network; however, even this is not particularly efficient. Thus, for portability reasons, SWAR programs should avoid the use of complex communication patterns.

2.3. Reductions

Reductions are operations in which the values stored in the fields of a register are combined into a single result value by successively applying an operation to an intermediate result and each of the field values. The final value may be stored in one or all of the fields of the result register.

For example, if 4x8 a partitioned register contains the values { 3, 4, 9, 18 }, we may store the result of a reduceAdd (3+4+9+18=34) in each of the fields to form the single result 34. In the following example, we perform an unsigned reduceAdd of the fields of the 4x8 partitioned register A containing the values { 4, 3, 2, 1 } to form the single result value 10:

	31	2	23	15		7	0
A:		4	3		2	Ĭ	1
B=A&mask:			3				1
C=A>>8:		0	4		3		2
D=C&mask:			4				2
E=B+D:			7				3
F=E>>16:			0				7
G=E+F:			7				10
G&mask:							10

This operation is performed recursively. First, the register is split into the even (multiples of two) and odd (non-multiples of two) fields using masking operations. Then, using an unpartitioned shift, the fields are aligned and added. The result is a register with half as many fields, but each is twice as large — conveniently ensuring that overflow will not occur, even if no spacers are used. This process is repeated until only one field remains.

2.4. Enable Masking

One of the distinguishing characteristics of SIMD computation is the ability to disable processing elements for portions of a computation. Unfortunately, SWAR hardware does not allow fields to be disabled per se; instead, a form of arithmetic nulling must be used. Consider the simple SIMD code fragment:

where (c)
$$a=b;$$

In this code, a, b, and c are all vectors of the same length. Where the corresponding field of c is true (non-zero), that field within a should be replaced with the value from the corresponding

field from b. Other elements of A should be unaltered. If it were possible to disable fields, the above statement could be executed very straightforwardly by disabling fields corresponding to 0 values within c and then simply having the enabled processing elements execute a=b;

Without the ability to disable fields, we will need to arithmetically nullify the undesired computations. The first step is to make the SIMD code fragment symmetric:

This looks strange, but accurately reflects the fact that the fields of a that are to be unaffected by the where must literally be actively read and pasted-together with the fields taken from b.

There are a variety of techniques that can be used to arithmetically merge the appropriate fields taken from a and b, but by far the cheapest is to use bitwise AND to mask the undesired field and then use bitwise OR to merge the masked results. Indeed, this is the approach used in the following example:

	31				23 15					7				0		
A:				4				3				2				1
В:				8				7				6				5
C:		0			27				148			0			0	
D=C>>4:		0		0		0		1		11		9		4		0
E=C/D:		0		0	1		11		11		13		4		0	
F=E>>2:	0	0	0	0	0	0	1	2	3	2	3	3	1	1	0	0
G=E F:	0	0	0	0	0	1	3	3	3	3	3	3	1	1	0	0
H=G>>1:	00	00	00	00	00	00	11	11	11	11	11	11	10	10	10	00
<i>I</i> =G/H:	00	00	00	00	00	0 1	11	11	11	11	11	11	11	11	10	00
J=I&LSBs:				0				1				1				0
K=MSBs-J:			1	28	127				127				128			
L=K^MSBs:		0			255			255			0					
M=B&L:				0				7				6				0
N=A&~L:				4				0				0				1
M/N:				4				7				6				1

First c must be converted into an appropriate mask, which has a field full of 0s where c was 0 and a field full of 1s where c was non-zero (in fact, although C logic true values ae set to 1, SWARC logic true values are set to -1 to simplify this use). A log-length sequence of unpartitioned shifts and bitwise OR operations is used to convert c into a usable mask, 1. This mask is then used to select the appropriate fields of c and c. In terms of scalar C code:

```
if (c) a=b; else a=a;
```

Essentially became:

```
1 = -(c != 0);

a = ((b & 1) | (a & ~1));
```

Although the SWAR code looks somewhat strange, this arithmetic masking is actually an old trick borrowed from various early SIMD machines. For example, the Thinking Machines CM-2 [TMC90] used this approach.

3. Compiler Optimizations For SWAR

Optimizations that have been devised for SIMD programming often apply in a natural way to SWAR programming. However, the unique features of SWAR execution also motivate a variety of new compiler technologies. Three such technologies are briefly discussed here: promotion of field sizes, SWAR value tracking, and enable masking optimizations.

3.1. Promotion Of Field Sizes

Just because a high-level language program stated that a value needed only *k* bits does not mean that *precisely k* bits worth of hardware must be used. For example, 16 bit values are handled very well by HP's PA-RISC MAX, but smaller sizes are not. Thus, a vector that was declared as containing 14-bit values will yield more efficient code sequences if the array's object type is promoted to 16-bit fields. This promotion is particularly favorable because the number of 16-bit fields that fit in a 64-bit register is precisely the same as the number of 14-bit fields that would fit — using 14-bit fields would not add any parallelism. Not all hardware-unsupported field sizes are inefficient; for example, 2-bit fields are not directly supported by any of the SWAR hardware, but for all the SWAR implementations described, the extra parallelism width makes operations on 2-bit fields far more effective than promoting these fields to 8-bits. In general, for each SWAR implementation, there are certain field sizes that are less efficient than a somewhat larger field size, and thus none of these inefficient field sizes should be directly used.

Notice that this promotion of field sizes, and even the use of spacer bits, can result in different data layouts for the same vector on different computers, and the user-specified field sizes do not imply any particular field layout. However, by using a separate SWAR module compiler that communicates only through C code interfaces that it generates, we can ensure that non-obvious machine-dependent layouts are not visible outside the SWAR modules.

3.2. SWAR Value Tracking

Nearly all traditional compiler optimizations are based on tracking which values are available where and when. For example, common subexpression elimination (CSE) depends on this analysis. The interesting question that SWAR technology raises is what constitutes the basic unit of SWAR data that we wish to track?

Fundamentally, one would like to track the values of all the fields within a SWAR register or vector. However, as we have discussed, the compiler's treatment of SWAR coding often results in code sequences that dynamically change the apparent partitioning, and this would have the effect of destroying the field-based value tracking information. Similarly, as discussed in sections 2.1.2 and 2.1.3, it is often desirable to use spacer bits between fields to allow ordinary instructions to function as partitioned operations, and these spacer bits are by definition not part of the fields they separate. To optimize spacer manipulations, we need to be able to track spacer values as well as field values. In fact, unpartitioned operations manipulate both field and spacer values alike, and operations such as shifts can transform one into the other.

To support aggressive SWAR optimizations, a new value tracking method is necessary. We suggest that symbolic tracking of values for arbitrary masked-bit-patterns within a register is appropriate. The following two subsections give brief examples of the benefits of such tracking.

3.2.1. Bitwise Value Tracking

It is not surprising that SWAR code often uses many bitwise masking (AND and OR) operations and unpartitioned shifts using constant-valued masks and shift distances. Consider a simple C code example:

```
x = (((x & 0x00ff) << 4) & 0xff00);
```

By tracking how arbitrary masked-bit-patterns are shifted in this sequence, this code can be converted to the equivalent, but simpler, form:

```
x = ((x << 4) \& 0x0f00);
```

In general, this type of tracking can merge multiple AND or OR operations through shifts, as well as merging shifts.

3.2.2. Simplification Of Spacer Manipulation

It is not unusual that the manipulations of spacer bits will become a significant fraction of all the instructions executed. Thus, any optimization technique that can reduce the frequency of spacer manipulations is highly desirable.

In preparation for a partitioned operation, the spacer bits may have to be set to particular values which depend on the operation. For example, if the operation is an addition, the spacer bits should be set to 0 in both operand registers. After a partitioned operation, a carry or borrow may alter these spacer bit values. Thus, it may be necessary to zero all the spacer bits to correctly isolate the fields.

Actually, most operations can alter the values of spacer bits. The interesting fact is that even though most polymorphic instructions, such as bitwise operations, can alter spacer bit values, they produce field values without being affected by the values of spacer bits. Thus, these instructions offer the opportunity to set spacer bits to the next desired value at no cost.

For example, consider computing e=((a+b)-(c+d)); using a SWAR representation employing spacer bits identified by the mask s:

```
e = (((((a \& ~s) + (b \& ~s)) \& ~s) | s) - ((((c \& ~s) + (d \& ~s)) \& ~s) \& ~s)) \& ~s;
```

We would expect conventional compiler optimizations to eliminate the redundant AND of the subtrahend with the one's complement of the spacer mask required by the field isolation stage of the addition and the normalization stage of the subtraction. This would save one operation, and change the calculation of the subtrahend from

```
((((c \& ~s) + (d \& ~s)) \& ~s) \& ~s)
```

To:

$$(((c \& ~s) + (d \& ~s)) \& ~s)$$

By using spacer value tracking, we may be able to make significantly larger reductions in the number of spacer manipulations needed. Suppose that the spacer bits for each of a, b, c, and d were already known by the compiler to be zeros. The normalizations for the additions would not be required, and the original code could be compacted to the following:

$$e = ((((a + b) \& ~s) | s) - (((c + d) \& ~s) \& ~s)) \& ~s;$$

Further inspection reveals that the isolation stages following the additions are unnecessary because they are immediately followed by the normalization stage of the subtraction, which overwrites the spacer values just written. This last observation reduces the original 12 operations to its final form, with just 6 operations (given that s and s are constants):

$$e = (((a + b) | s) | ((c + d) & ~s)) & ~s;$$

It is also interesting to note that by the same analysis, computing e=((a+b)-((c+d)&M)); where M is a constant-valued mask, could be accomplished in the same number of instructions:

$$e = (((a + b) | s) | ((c + d) & (M & ~s)) & ~s;$$

This is because $(M \& ^s)$ is also a constant. In general, the bitwise operations and unpartitioned shifts discussed in the previous section can be optimized at the same time as the spacer manipulation.

3.3. Enable Masking Optimizations

Because SWAR hardware does not allow fields to be disabled, there is a significant cost associated with the arithmetic nulling of undesired field computations. However, this cost need not be incurred if the compiler's static analysis can prove that *all* fields are active, or if the compiler can generate code that allows all fields to be active later correcting the should-have-been-inactive field values.

For "virtualized" processing elements, a single vector may span the fields of multiple words (registers). Thus, enable masking would have to be performed on each of the words. This allows the compiler to generate three different versions of the word-wide SWAR code, and to select among them by examining each word worth of enable mask:

- If the enable mask word is entirely enabled, then no masking is done; the corresponding word of the result is directly computed.
- If the enable mask word is entirely disabled, then the corresponding word of the result is copied from the original value no computation is done.
- If the enable mask word is partially enabled, then the usual masking is used in computing the result.

4. A SWAR High-Level Language And Its Compiler

Although compiler analysis could attempt to recognize the number of bits precision needed for each object, and well-known compiler vectorization transformations could be applied to a sequential language, these are secondary concerns. The first step is to design, build, and optimize a portable language and compilers that allow the programmer to explicitly specify both object precision and SIMD-style parallelism. For this purpose, a C-like language called SWARC (pronounced "swh-are-see") was created. The complete syntax of SWARC is beyond the scope of this paper; here we consider only the most important extensions beyond C's semantics and how the Scc compiler produces efficient SWAR code using the MMX instructions.

4.1. The SWARC Language

The key concept in SWAR is the use of operations on word-sized values to perform SIMD-parallel operations on fields within a word. The field sizes and alignments are highly machine dependent, so, for portability, the programmer's specification should not *require* particular layouts, but should provide *minimal constraints* on how a particular layout can be selected by the compiler. To interface to C as a module compiler, SWARC also needs to be able to integrate well with ordinary C code and data structures. Thus, although SWARC is based on C, these design issues lead to adjustments in the type system, operators, and general structure of the language.

4.1.1. Type System

A SWARC data type specifies a first-class array type whose elements are either ordinary C-layout objects or are SWAR-layout integer or floating-point fields with a *specified minimum precision*. For example, the C declaration <code>char c[14]</code>; is equivalent to the SWARC C-layout declaration <code>char[14] c;</code>. The SWAR-layout declaration syntax is similar to the bit-field specification used in C <code>struct</code> declarations; <code>char:7[14] d;</code> declares d to be an array with 14 visible elements, each of which is an integer field with *at least* 7 bits precision. Using the precision specifier without an integer precision is equivalent to specifying the native precision for the equivalent C-layout type, e.g., <code>float:[5]</code> is equivalent to <code>float:32[5]</code>.

In C, operations on all integer types use modular arithmetic. However, operations on natural data types often require saturation arithmetic and such operations are well-supported by many SWAR targets. SWARC allows either modular or saturation arithmetic to be specified as an attribute of a type. Even if the precision actually used by the compiler for an object is greater that the minimum precision specified by the user, results saturate at the specified number of bits.

Because of these type extensions, the SWARC type coercion rules are somewhat different from those of C:

- 1. Width. An expression in which one operand is a scalar (has width one), and the other operand has a greater width, yields a result of the greater width. If neither object is a scalar and the widths do not match, a warning is generated and the wider object is truncated.
- 2. *Precision*. An expression in which one operand has a C-layout and the other has a SWAR-layout results in the SWAR-layout even if precision is reduced; explicit precision overrides implicit precision. Otherwise, an expression with mixed precision operands yields a result with the higher precision.
- 3. Sign. Mixed signed and unsigned operands yield a signed result.
- 4. Arithmetic Mode. Mixed modular and saturation operands yield a saturation result.

4.1.2. Operators

Most C operators have an obvious and useful meaning as element-by-element operations on SWARC first-class arrays. In addition, SWARC adopts the C^* [TMC90a][HaL91] concept of using assignment operators to perform reductions when storing an array into a scalar or when using the assignment operator as a unary prefix, although SWARC only allows such use for associative operators. For example, +=c would yield the scalar sum of the elements within c.

SWARC also adds a few operators. The widthof and precision of operators can be applied to an expression to return constant values that are, respectively, the declared number of elements or bits of precision. There are a variety of unary "communication" operations; for example, c[<<% 3] returns an array value corresponding to a circular shift by three of the elements of c. Minimum (?<), maximum (?>), and average (+/) binary operators are

particularly useful. These binary operators also are allowed to be used as assignment and reduction operators.

The fact that SWARC is essentially a SIMD language also yields a minor semantic adjustment: the "short-circuit" evaluation of C operators ? ... :, $|\cdot|$, and && using jumps is required for C, but is optional for SWARC. Unlike C and C*, SWARC also allows $|\cdot|$ = and &&= to be used for assignments and reductions (as reductions, these are the well-known SIMD ANY and ALL operators).

4.1.3. Functions, Control Flow, and General Structure

Aside from the differences in typing, SWARC function definitions look identical to those of ANSI C. However, to minimize copy overhead, SWARC uses call by address, whereas C uses call by value.

The control flow constructs of SWARC are a superset of those in C, with a superset of the semantics. For example, if, while, and for statements with array-valued conditions work much as those of MPL [Mas91]: an enable mask is computed and the body is executed only if ANY element is enabled. There is also a where ... elsewhere construct which behaves like if, but without the ANY test; the everywhere statement provides a way for a nested construct to temporarily enable all elements. Array operations lexically within an enable masking construct respect the mask; called functions, however, do not. To ensure compatibility with ordinary C code, every function definition begins with everywhere implicitly enabled.

The C module language design of SWARC also is reflected in the fact that SWARC allows arbitrary C code to be embedded within SWARC code in the syntactic role of either statements or global declarations. Embedded C code is surrounded by \${ ... \$}; although # C-preprocessor directives are used to preprocess SWARC code, the embedded C code can use \$ to specify later (nested) preprocessing to be applied to the final C code. For example, \$include <file> would include file when compiling the C code, whereas #include <file> would include file in the SWARC compile.

4.2. Scc, The SWARC Compiler

Scc is the first SWARC compiler. Currently, it consists of approximately 50,000 lines of C code, written primarily as a prototype testbed for SWAR compiler technology. However, Scc is expected to evolve into a full public domain release stable enough for production work. Although it will eventually support many different SWAR targets, at this writing it targets 32-bit C (performing SWAR operations using ordinary C code) and, using nearly cycle-accurate models, seven different MMX variants: AMD K6, AMD K6-2, Cyrix 6x86MX, Cyrix MII, Intel Pentium II, Intel Pentium MMX, and generic MMX (an averaged target derived from the others).

4.2.1. Approach

Although some portions of Scc are complex and unusual, the basic approach is quite conventional.

An LL(2) parser, created using PCCTS (the Purdue Compiler Construction Tool Set, see network newsgroup <code>comp.compilers.tools.pccts</code>), is used to recognize the SWARC language. As the program is parsed, an arbitrary-precision arbitrary-length vector intermediate representation is created, type coerced, and simple optimizations (constant folding, etc.) are applied.

When the end of a function is reached, the vector intermediate code for the function is passed to the *fragmenter*. This translates the vector code into machine-specific optimized RISC-like operation DAGs operating on fixed-size word values. No vector temporaries are generated; the only temporaries generated are the register results of DAG operations on words. At this point the actual precisions and SWAR-layouts are determined and internally fixed, selecting from 1, 2, 4, 8, 16, or 32-bit densely-packed SWAR field layouts. Optimizations include the usual basic block techniques (common subexpression elimination, value propagation, dead code elimination, etc.), a variety of machine-specific translations, and a growing number of SWAR-specific new optimizations (such as the bitwise value tracking discussed earlier in this paper).

For MMX-based targets, there are only eight 64-bit registers, so register allocation is critical. The fact that MMX instructions allow one memory operand (i.e., rather than referencing a register that was explicitly loaded) can be used to reduce register pressure, but this interferes with pipeline scheduling. Machine-specific pipeline constraints are very important, especially in processors that have two MMX pipelines. The relatively small speedups offered by SWAR technology could be destroyed by either bad scheduling or register spills.

Our solution *never* spills a register and often generates provably optimal schedules. First, a heuristic prescheduler finds a schedule order with a low maxlive by working backward from word stores to partition the schedule. Using this initial schedule, Scc applies a modified exhaustive search combined scheduler and register-allocator based on schedule permutation [NiD90]. If allowed to run to completion, the resulting schedule and register allocation is optimal for the detailed pipeline model used; by default, the search stops with the first valid schedule. A maximum search time can be given on the Scc command line; typical blocks of up to hundreds of instructions can be exhaustively searched in under a second. In the very rare case that the scheduler cannot find a spill-free schedule, it simply reports an internal scheduling failure; future Scc versions will insert spills only in this case.

4.2.2. Compiled Code

Space does not permit examples of all the code generation mechanisms, but a simple example gives the general flavor. Consider the following one-line SWARC function:

```
void f(char[14] a, char:[14] b) { a += b; }
```

The code output by Scc is C code with macros wrapping the inline native assembly code for each MMX instruction. When this code is compiled using gcc, the MMX assembly language code incurs no overhead from the macro inlining. Targeting generic MMX for the above example, the code is:

```
#include "Sc.h"
void f(char *a,
mmx_t *b)
    extern mmx_t mmx_cpool[];
    register mmx_t *_cpool = &(mmx_cpool[0]);
        movq_m2r(*(((mmx_t *) a) + 0), mm0);
        movq_m2r(*(((mmx_t *) a) + 1), mm1);
        paddb_m2r(*(((mmx_t *) b) + 0), mm0);
        movq_r2r(mm1, mm2);
        paddb_m2r(*(((mmx_t *) b) + 1), mm1);
        movq_r2m(mm0, *(((mmx_t *) a) + 0));
        pand_m2r(*(_cpool + 2), mm1);
        pand_m2r(*(\_cpool + 3), mm2);
        por_r2r(mm2, mm1);
        movq_r2m(mm1, *(((mmx_t *) a) + 1));
    }
_return:
    emms();
/* MMX constant pool */
mmx_t mmx_cpool[] = {
   0 */ 0x00000000000000LL,
               0xffffffffffffffLL,
     1 */
   1 */ 0xffffffffffffffLL,

2 */ 0x0000ffffffffffffLL,

3 */ 0xffff000000000000LL
/*
```

At first, the code seems more complex than necessary; however, this complexity comes from the fact that the code stores into the C-layout variable a and 14 characters does not fill two words. Thus, although the addition of the first words of a and b is trivial, before the result from adding the second two words can be stored it is necessary that the 14th and 15th character positions be set to the values originally in those memory locations. This masking, which would not be needed if we were storing into a SWAR-layout object, is done by the pand and por operations. The constant pool is needed because the MMX instruction set does not provide an immediate load; in fact, the _cpool pointer is used to ensure that the constant pool addressing within the MMX instructions will be able to use a short offset (because a long address would disrupt the instruction decode pipeline).

An interesting note is that, with better bitwise value tracking, this good code could be made better: addition of 0 preserves the original value, thus, a single pand with the second word of b before adding could have had the same effect as the masking that was applied later using three instructions.

5. Conclusion

The latest processors from AMD, Cyrix, DEC, HP, Intel, MIPS, Motorola, and Sun have proven that hardware support for SWAR, SIMD Within A Register, is easily and efficiently added to a conventional microprocessor. However, more than year after most of these processors became available, there are still no general-purpose high-level languages that allow them to be programmed using a portable SIMD model. Essentially by design, these SWAR hardware implementations are tuned for hand-coding of specific algorithms, with sparse coverage of SWAR functionality driven by the need to minimize disturbance of the base processor instruction set and architecture.

In this paper, we have shown that these flaws can be largely overcome by the combination of clever coding sequences and a few new types of compiler analysis and optimization. Only the lack of random inter-field communication is an unsolvable problem, but these operations have been avoided in SIMD algorithms because many traditional SIMD machines also lacked the necessary hardware support.

The SWARC module language and Scc prototype compiler, both briefly overviewed here, are good first steps toward making full SIMD-style programming models effective for commodity microprocessors. It is our intention that Scc will mature into a public domain production-quality tool that can be used not only by programmers, but also as a target for higher-level parallelizing compilers. This is especially important given the introduction of high-performance SWAR floating point in AMD's K6-2 3DNow! and other processors.

Further information about the SWARC language and Scc compiler is available on the WWW at http://shay.ecn.purdue.edu/~swar/Index.html.

References

[AMD97] Advanced Micro Devices, Inc., AMD-K6 Processor Multimedia Extensions, Advanced Micro Devices, Inc., Sunnyvale, California, March 1997.

[AMD98] Advanced Micro Devices, Inc., 3DNow! Technology Manual, Advanced Micro Devices, Inc., http://http.amd.com/K6/k6docs/pdf/21928c.pdf, May 1998.

[Cyr96] Cyrix Corporation, *Multimedia Instruction Set Extensions for a Sixth-Generation x86 Processor*, Cyrix Corporation, ftp://ftp.cyrix.com/developr/hc-mmx4.pdf, August 1996.

[Cyr98] Cyrix Corporation, *Application Note 108: Cyrix Extended MMX Instruction Set*, Cyrix Corporation, http://http.national.com/cyrix/mii/108ap.pdf, 1998.

[DEC96] Digital Equipment Corporation, *Alpha Architecture Handbook, Version 3*, Digital Equipment Corporation, Maynard, Massachusetts, October 1996.

[HaL91] P. J. Hatcher, A. J. Lapadula, R. R. Jones, M. J. Quinn, and R. J. Anderson, "A Production-Quality C* Compiler for Hypercube Multicomputers" *Third ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming*, Williamsburg, Virginia, April 1991, pp. 73-82.

[CeT97] Centaur Technology, Inc., *IDT WinChip C6 Processor Data Sheet*, version 1.00, Integrated Device Technology Inc., http://www.winchip.com/pdf/C6DS.pdf, 1997.

[Int96] Intel Corporation, *Intel Architecture MMX Technology: Programmer's Reference Manual*, Intel Corporation, http://developer.intel.com/drg/mmx/Manuals/prm/prm_covr.htm, March 1996.

[Lee95] Ruby B. Lee, Accelerating Multimedia with Enhanced Microprocessors, IEEE Micro, 15(2):22-32, April 1995

[LeH96] Ruby Lee and Jerry Huck, *HP Technical Computing* — 64-bit and Multimedia in the PA-RISC 2.0 Architecture, Hewlett-Packard Company, http://hpcc997.external.hp.com:80/wsg/strategies/pa2go3/pa2go3.html, June 1996.

[Mas91] MasPar Computer Corporation, MasPar Programming Language (ANSI C compatible MPL) Reference Manual, Software Version 2.2, Document Number 9302-0001, Sunnyvale, California, November 1991.

[Mot98] Motorola Inc., *AltiVec Technology Programming Environments Manual, Rev. 0.2*, Motorola Inc., http://www.mot.com/SPS/PowerPC/teksupport/teklibrary/manuals/altivec_pem.pdf, May 1998.

[NiD90] A. Nisar and H.G. Dietz, "Optimal Code Scheduling for Multiple Pipeline Processors," 1990 International Conference on Parallel Processing, vol. II, Saint Charles, Illinois, pp. 61-64, August 1990.

[SGI96] Silicon Graphics, Inc., MIPS Digital Media Extension, Silicon Graphics, Inc., http://www.sgi.com/MIPS/arch/ISA5/, 1996.

[Sun97] Sun Microsystems, Inc., *The VIS Instruction Set, Sun Microelectronics*, Sun Microsystems Corporation, http://www.sun.com/sparc/vis/index.html, April 1997.

[TMC90] Thinking Machines Corporation, Connection Machine Model CM-2 Technical Summary, version 6.0, Thinking Machines Corporation, Cambridge, Massachusetts, November 1990.

[TMC90a] Thinking Machines Corporation, *C* Programming Guide*, Thinking Machines Corporation, Cambridge, Massachusetts, November 1990.