

CS:APP3e Web Aside OPT:SIMD: Achieving Greater Parallelism with SIMD Instructions*

Randal E. Bryant
David R. O'Hallaron

December 29, 2014

Notice

The material in this document is supplementary material to the book Computer Systems, A Programmer's Perspective, Third Edition, by Randal E. Bryant and David R. O'Hallaron, published by Prentice-Hall and copyrighted 2016. In this document, all references beginning with "CS:APP3e " are to this book. More information about the book is available at csapp.cs.cmu.edu.

This document is being made available to the public, subject to copyright provisions. You are free to copy and distribute it, but you must give attribution for any use of this material.

1 Introduction

As described in CS:APP3e Section 3.1, Intel introduced the SSE instructions in 1999, where SSE is the acronym for "streaming extensions," and, in turn, SIMD (pronounced "sim-dee") is the acronym for "single-instruction, multiple-data." implement the AVX multimedia instructions. In addition to providing a superset of the SSE instructions Our presentation is based on AVX2, the second version of AVX, introduced with the Core i7 Haswell processor in 2013. GCC will generate AVX2 code when given the command-line parameter `-mavx2`. We will show how the SIMD instructions can be used to implement *vector* computations, where multiple data values are computed with a single operation.

The idea behind the SIMD execution model in AVX2 is that each 32-byte YMM register can hold multiple values. In our examples, we will consider the cases where they can hold either eight integer or single-precision values, or four double-precision values. AVX2 instructions can then perform vector operations on these registers, such as adding or multiplying eight or four sets of values in parallel. For example, if YMM register `%ymm0` contains eight single-precision floating-point numbers, which we denote a_0, \dots, a_7 , and `%rcx` contains the memory address of a sequence of eight single-precision floating-point numbers, which we denote b_0, \dots, b_7 then the instruction

*Copyright © 2015, R. E. Bryant, D. R. O'Hallaron. All rights reserved.

```
vmulps  (%rcx), %ymm0, %ymm1
```

will read eight values from memory and eight from AVX register `%ymm0`. It will perform eight multiplications in parallel, computing $c_i \leftarrow a_i \cdot b_i$, for $0 \leq i \leq 7$ and store these results in AVX register `%ymm1`. We see that a single instruction is able to generate a computation over multiple data values, hence the term “SIMD.” This multiplication is an example of what we will refer to as *vector code*. We will refer to code that operates only on one value at a time as *scalar code*.

GCC supports extensions to the C language that let programmers express a program in terms of vector operations that can be compiled into the SIMD instructions of AVX2 [1]. This coding style is preferable to writing code directly in assembly language, since GCC can also generate code for the SIMD instructions found on other processors. Writing in C also has the advantage that GCC will generate scalar code for machines that do not support vector instructions. We will describe how to write code using the GCC support for vector operations, using our combining functions as examples. The basic strategy is to define a vector data type `vec_t` that holds either eight 4-byte values or four 8-byte values. If we have two such vectors `va` and `vb`, then the expression `va * vb` causes a SIMD multiplication of the vector elements.

2 Declarations

Our first step is to declare the vector data type. Since we are trying to make the same code work for base data types `int`, `long`, `float`, and `double`, we will use a combination of `typedef` declarations and constant definitions to make the code more general. As with our earlier versions, we assume that the base data type has been declared as type `data_t`.

We define `VBYTES` to be the number of bytes in a vector. For AVX2, this is defined to be 32, but we would like to keep this value parameterized to easily adapt the code for other machines. We then defined `VSIZE` to be the number of elements in each vector:

```
1 /* Number of bytes in a vector */
2 #define VBYTES 32
3
4 /* Number of elements in a vector */
5 #define VSIZE VBYTES/sizeof(data_t)
```

We are now ready to define the vector data type. This involves a notation that is idiosyncratic to GCC:

```
/* Vector data type */
typedef data_t vec_t __attribute__((vector_size(VBYTES)));
```

This declaration states that data type `vec_t` is vector, where the elements are of type `data_t`, and the vector consists of `VBYTES` bytes.

We now need some means of accessing the elements of a vector. Rather than introducing additional notation, we can make use of the `union` declaration to overlay vector and array data types:

```
typedef union {
```

```

    vec_t v;
    data_t d[VSIZE];
} pack_t;

```

The following shows a simple example of computing the inner product of two SIMD vectors:

```

1 /* Compute inner product of SSE vector */
2 data_t innerv(vec_t av, vec_t bv) {
3     pack_t xfer;
4     long int i;
5     vec_t pv = av * bv;
6     data_t result = 0;
7     xfer.v = pv;
8     for (i = 0; i < VSIZE; i++)
9         result += xfer.d[i];
10    return result;
11 }

```

In this code, the multiplication operation on line 5 multiplies the corresponding elements of vectors `av` and `bv`. The code on lines 7–10 then accesses and sums the elements of the product vector. Using the variable `xfer`, of type `pack_t`, the code provides access to each element i of the vector with the expression `xfer.d[i]`.

3 Alignment Requirement

Some of the AVX2 instructions impose a very strict alignment requirement on memory operands. They require that any data being read from memory into an YMM register, or written from an YMM register to memory, satisfy a 32-byte alignment. An instruction that attempts to read or write unaligned data can cause a segmentation fault, indicating an invalid memory reference. This alignment requirement will factor into how we write programs that make use of AVX2 instructions. (There are AVX2 instructions that can access unaligned data, but early implementations were not very efficient, and so GCC does not currently generate code that uses them.)

4 Implementation of Combining Function

Figure 1 shows the code for a combining function that makes use of the SIMD operations. The overall idea is to set up a vector variable `accum` that accumulates either eight (data types `int` and `float`) or four (data type `double`) values in parallel.

The code starts by initializing the accumulators to the identity element (lines 11–13), using the `pack_t` data type to set the individual elements of a vector.

To satisfy the alignment requirement, we may need to accumulate several vector elements using scalar operations until the remaining data vector `data` has an address that is a multiple of `VBYTES`. The code for this is shown in lines 16–19. Observe we cast pointer `data` to data type `size_t` so that we can test

```

1 void simd_vl_combine(vec_ptr v, data_t *dest)
2 {
3     long i;
4     pack_t xfer;
5     vec_t accum;
6     data_t *data = get_vec_start(v);
7     int cnt = vec_length(v);
8     data_t result = IDENT;
9
10    /* Initialize accum entries to IDENT */
11    for (i = 0; i < VSIZE; i++)
12        xfer.d[i] = IDENT;
13    accum = xfer.v;
14
15    /* Single step until have memory alignment */
16    while (((size_t) data) % VBYTES) != 0 && cnt) {
17        result = result OP *data++;
18        cnt--;
19    }
20
21    /* Step through data with VSIZE-way parallelism */
22    while (cnt >= VSIZE) {
23        vec_t chunk = *((vec_t *) data);
24        accum = accum OP chunk;
25        data += VSIZE;
26        cnt -= VSIZE;
27    }
28
29    /* Single-step through remaining elements */
30    while (cnt) {
31        result = result OP *data++;
32        cnt--;
33    }
34
35    /* Combine elements of accumulator vector */
36    xfer.v = accum;
37    for (i = 0; i < VSIZE; i++)
38        result = result OP xfer.d[i];
39
40    /* Store result */
41    *dest = result;
42 }

```

Figure 1: **Combining function using SIMD operations.** The vector operations cause multiple (4 or 8) values to be accumulated in parallel in variable `accum`.

whether it is a multiple of `VBYTES`. We must also keep track of the number of remaining elements `cnt`, and consider the case where `cnt` is smaller than the number of elements in a single vector.

Lines 22–27 show the main loop for the function. We see here the use of casting to create a pointer to a vector having the same address as the pointer to the data. Dereferencing this pointer then retrieves an entire vector of data from memory, defined here by vector variable `chunk`. The statement `accum = accum OP chunk` then combines the vector values read from memory with the values in the parallel accumulators.

In the event that the main loop terminates before all values have been accumulated, we have another loop to single step through the remaining elements (lines 30–33.)

We then reference the accumulators through a `union` and combine them to accumulate the final result (lines 36–38.)

5 Analysis

The following table shows our results for the code of Figure 1, compared to our best method using only scalar operations:

Method	Integer				Floating point			
	int		long		float		double	
	+	*	+	*	+	*	+	*
Scalar 10×10	0.54	1.01	0.55	1.00	1.01	0.51	1.01	0.52
Scalar throughput bound	0.50	0.50	1.00	1.00	1.00	1.00	0.50	0.50
Vector 1×1	0.16	1.25	2.14	4.51	0.37	0.63	0.75	1.26
Vector throughput bound	0.06	0.12	0.12	–	0.12	0.06	0.25	0.12

Note that we now list the performance for the different integer sizes (`int` and `long`) and the different floating-point sizes (`float` and `double`) separately. Although these two cases have identical performance for scalar code, they have different performance with vector code, since one achieves 8-way parallelism, while the other just 4-way. We see that the resulting performance is a bit mixed. For integer addition and both single-precision addition and multiplication, we have broken the throughput barrier that we have seen for scalar implementations. On the other hand, other operations did not perform as well as we achieved using scalar code.

Fortunately, we can combine our earlier techniques for further enhancing parallelism either by expanding the number of accumulators (see Problem 1) or by using reassociation (see Problem 2), yielding the following performance:

Method	Integer				Floating point			
	int		long		float		double	
	+	*	+	*	+	*	+	*
Scalar 10×10	0.54	1.01	0.55	1.00	1.01	0.51	1.01	0.52
Scalar throughput bound	0.50	0.50	1.00	1.00	1.00	1.00	0.50	0.50
Vector 8×8	0.05	0.24	0.13	1.51	0.12	0.08	0.25	0.16
Vector $8 \times 1a$	0.06	0.25	0.25	1.87	0.51	0.27	1.04	0.52
Vector throughput bound	0.06	0.12	0.12	–	0.12	0.06	0.25	0.12

We can see that the vector code achieves almost an eightfold improvement on the four 32-bit cases, and a fourfold improvement on three of the four 64-bit cases. Only the long integer multiplication code does not perform well when we attempt to express it in vector code. The AVX instruction set does not include one to do parallel multiplication of 64-bit integers, and so GCC cannot generate vector code for this case. Using vector instructions creates a new throughput bound for the combining operations. These are eight times lower for 32-bit operations and four times lower for 64-bit operations than the scalar limits. Our code comes close to achieving these bounds for several combinations of data type and operation.

6 Problems

Practice Problem 1:

Write vector code for the combining function that maintains four different sets of accumulators, each accumulating either eight or four values, depending on the data type.

Practice Problem 2:

Write vector code for the combining function that reads four 32-byte chunks from memory on each iteration, and then uses reassociation to increase the number of combining operations that can be performed in parallel.

Practice Problem 3:

Write a SIMD version of the inner-product computation described in CS:APP3e Problem 5.13, using a single vector variable to accumulate multiple sums in parallel. You cannot assume that the argument vectors satisfy a 32-byte alignment. You can assume, however, that any degree of misalignment will be the same for both. In other words, for pointer p , the expression $((\text{long})\ p) \% 32$ will yield the same value when p is `udata` as it will when p is `vdata`.

Our implementation of this function achieves a CPE of 0.38 for single-precision data, and 0.75 for double-precision data.

Practice Problem 4:

Extend your code for Problem 3 to accumulate sums in four vectors. Our implementation of this function achieves a CPE of around 0.14 for single-precision data, and 0.29 for double-precision data, nearly reaching the throughput limit imposed by the add unit.

Practice Problem 5:

Write a SIMD version of the polynomial evaluation described in CS:APP3e Problem 5.5, using a single vector variable to accumulate multiple sums in parallel. Your code must work correctly regardless of the alignment of argument *a*. Our implementation of this function achieves a CPE of 1.44, exceeding the performance of the best scalar implementation.

Practice Problem 6:

Use the various tricks you have learned: vector code, multiple accumulators, reassociation, and Horner's method to write the fastest polynomial evaluation function you can. Our code achieves a CPE of 0.35.

Practice Problem Solutions**Problem 1 Solution: [Pg. 6]**

This code combines the style we have seen for parallel accumulation with vector code. Here is the main loop for the function:

```
/* 4 * VSIZE x 4 * VSIZE unrolling */
while (cnt >= 4*VSIZE) {
    vec_t chunk0 = *((vec_t *) data);
    vec_t chunk1 = *((vec_t *) (data+VSIZE));
    vec_t chunk2 = *((vec_t *) (data+2*VSIZE));
    vec_t chunk3 = *((vec_t *) (data+3*VSIZE));
    accum0 = accum0 OP chunk0;
    accum1 = accum1 OP chunk1;
    accum2 = accum2 OP chunk2;
    accum3 = accum3 OP chunk3;
    data += 4*VSIZE;
    cnt -= 4*VSIZE;
}
```

The following code then combines the results for all of the accumulators:

```
/* Combine into single accumulator */
xfer.v = (accum0 OP accum1) OP (accum2 OP accum3);

/* Combine results from accumulators within vector */
for (i = 0; i < VSIZE; i++)
    result = result OP xfer.d[i];
```

Problem 2 Solution: [Pg. 6]

This version only requires modifying the main loop:

```
while (cnt >= 4*VSIZE) {
    vec_t chunk0 = *((vec_t *) data);
```

```

    vec_t chunk1 = *((vec_t *) (data+VSIZE));
    vec_t chunk2 = *((vec_t *) (data+2*VSIZE));
    vec_t chunk3 = *((vec_t *) (data+3*VSIZE));
    accum = accum OP
        ((chunk0 OP chunk1) OP (chunk2 OP chunk3));
    data += 4*VSIZE;
    cnt -= 4*VSIZE;
}

```

As this example shows, we can enhance parallelism by reassociating the vector operations, just as we did for scalar operations.

Problem 3 Solution: [Pg. 6]

This code is a direct adaptation of the SIMD version of the combining function:

```

1  /* Vectors must satisfy alignment requirement */
2  void inner_simd_v1(vec_ptr u, vec_ptr v, data_t *dest)
3  {
4      long i;
5      pack_t xfer;
6      vec_t accum;
7      data_t *udata = get_vec_start(u);
8      data_t *vdata = get_vec_start(v);
9
10     long cnt = vec_length(v);
11     data_t result = 0;
12
13     /* Initialize accum entries to 0 */
14     for (i = 0; i < VSIZE; i++)
15         xfer.d[i] = 0;
16     accum = xfer.v;
17
18     /* Single step until have memory alignment */
19     while (((long) udata) % VBYTES && cnt) {
20         result += *udata++ * *vdata++;
21         cnt--;
22     }
23
24     /* Step through data with VSIZE-way parallelism */
25     while (cnt >= VSIZE) {
26         vec_t uchunk = *((vec_t *) udata);
27         vec_t vchunk = *((vec_t *) vdata);
28
29         accum = accum + (uchunk * vchunk);
30         udata += VSIZE;
31         vdata += VSIZE;
32         cnt -= VSIZE;
33     }
34

```

```

35     /* Single-step through remaining elements */
36     while (cnt) {
37         result += *udata++ * *vdata++;
38         cnt--;
39     }
40
41     /* Combine elements of accumulator vector */
42     xfer.v = accum;
43     for (i = 0; i < VSIZE; i++)
44         result += xfer.d[i];
45
46     /* Store result */
47     *dest = result;
48 }

```

Problem 4 Solution: [Pg. 6]

The following shows the inner loop for this function

```

/* Step through data with 4*VSIZE-way parallelism */
while (cnt >= 4*VSIZE) {
    vec_t uchunk = *((vec_t *) udata);
    vec_t vchunk = *((vec_t *) vdata);
    accum0 = accum0 + (uchunk * vchunk);
    udata += VSIZE; vdata += VSIZE;

    uchunk = *((vec_t *) udata);
    vchunk = *((vec_t *) vdata);
    accum1 = accum1 + (uchunk * vchunk);
    udata += VSIZE; vdata += VSIZE;

    uchunk = *((vec_t *) udata);
    vchunk = *((vec_t *) vdata);
    accum2 = accum2 + (uchunk * vchunk);
    udata += VSIZE; vdata += VSIZE;

    uchunk = *((vec_t *) udata);
    vchunk = *((vec_t *) vdata);
    accum3 = accum3 + (uchunk * vchunk);
    udata += VSIZE; vdata += VSIZE;

    cnt -= 4*VSIZE;
}

```

Problem 5 Solution: [Pg. 7]

Our code accumulates multiple sums in parallel, using a vector `xpwrv` with elements x^i, x^{i+1}, \dots when executing the loop where pointer `a` is the address of coefficient a_i .

```

1 double poly_simd_v1(double a[], double x, long degree)

```

```

2 {
3     long i;
4     pack_t xfer;
5     vec_t accum;
6     long cnt = degree+1;
7     double result = 0;
8     double xpwr = 1.0; /* Various powers of x */
9
10    vec_t xvv; /* Vector of x^{VSIZE} */
11    vec_t xpwrv; /* Vector of increasing powers of x */
12
13    /* Initialize accum entries to 0, and compute x^{VSIZE} */
14    xpwr = 1.0;
15    for (i = 0; i < VSIZE; i++) {
16        xfer.d[i] = 0;
17        xpwr *= x;
18    }
19    accum = xfer.v;
20
21    /* Create a vector of all x^{VSIZE} */
22    for (i = 0; i < VSIZE; i++)
23        xfer.d[i] = xpwr;
24    xvv = xfer.v;
25
26    xpwr = 1;
27    /* Single step until have memory alignment */
28    while (((size_t) a) % VBYTES && cnt) {
29        result += *a++ * xpwr;
30        xpwr *= x;
31        cnt--;
32    }
33
34    /* Create vector with values xpwr, xpwr*v, ... */
35    for (i = 0; i < VSIZE; i++) {
36        xfer.d[i] = xpwr;
37        xpwr *= x;
38    }
39    xpwrv = xfer.v;
40
41    /* Main loop. Accumulate sums in parallel */
42    while (cnt >= VSIZE) {
43        vec_t chunk = *((vec_t *) a);
44        accum += chunk * xpwrv;
45        xpwrv *= xvv;
46        a += VSIZE;
47        cnt -= VSIZE;
48    }
49
50    /* Extract accumulated values */
51    xfer.v = accum;

```

```
52     xpwr = 1.0;
53     for (i = 0; i < VSIZE; i++)
54         result += xfer.d[i];
55
56     /* Get highest power of x */
57     xfer.v = xpwr;
58     xpwr = xfer.d[0];
59     /* Single step remaining elements */
60     while (cnt) {
61         result += *a++ * xpwr;
62         xpwr *= x;
63         cnt--;
64     }
65
66     return result;
67 }
```

Problem 6 Solution: [Pg. 7]

Solution available on Instructor's portion of CS:APP website.

References

[1] *GCC Online Documentation*. Available at <http://gcc.gnu.org/>.