

Intrinsics for SIMD

Part 2

Satish Puri

Evaluate Polynomial

$$3x^9 + x^7 + 12x^6 + 8x^5 + 7x^4 + x^2 - 9x - 6$$

| Array Index | 9 | 8 | 7 | 6 | 5 | 4 | 2 | 1 | 0 |
|-------------|-----|-----|-----|----|----|----|---|----|----|
| coefficient | 3 | 0 | 1 | 12 | 8 | 7 | 1 | -9 | -6 |
| exponent | 512 | 256 | 128 | 64 | 32 | 16 | 4 | 2 | 1 |
| Term | | | | | | | | | |

$x = 2$

```
double total = 0;
for(int i = 0; i < TERMS; i++)
{
    Term t = coefficient[i] * exponent[i]
    Add t to total
}
```

Using intrinsics for Polynomial Evaluation

- `double coeff[MAX] = {1,2,4, 8, 5, ..., 88};`
- `double expon[MAX] = {512, 256, 128, 64, ..., 4, 2,1 };`
- `double term[MAX]`
- `For(int t = 0; t< MAX; t++)`
 - `{`
 - `term[t] = coeff[t] * expon[t];`
 - `}`

Array program

```
void vecMulAdd(int n, double *A, double *B, double *C)
{
    int i;
    for(i = 0; i<n; i++)
    {
        C[i] = C[i] + A[i] * B[i];
    }
}
```

A 1 2 3 4 5 6 7 8

B 0 1 2 3 6 7 8 9

C 0 0 0 0 0 0 0 0 0

0 2 6 12 30 42 56 72

Translate the code using SIMD Intrinsics

```
void vecMulAdd(int n, double *A, double *B, double *C)
{
    int i;
    for(i = 0; i<n; i++)
    {
        C[i] = C[i] + A[i] * B[i];
    }
}
```

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| A | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| B | 0 | 1 | 2 | 3 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| C | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|

| | | | | | | | |
|---|---|---|----|----|----|----|----|
| 0 | 2 | 6 | 12 | 30 | 42 | 56 | 72 |
|---|---|---|----|----|----|----|----|

Memory alignment requirements

For doubles, 32 bytes aligned
memory address starts at modulo 32

`_mm_malloc(amount of memory, alignment)`
ensures that data is aligned to a 32 byte boundary.

- An instruction that attempts to read or write unaligned data can cause a segmentation fault, indicating an invalid memory reference.
- Poor performance if not aligned.

Alignment question 1

Let us assume that **arr** array is stored in a block of memory with starting location $(3200)_{10}$

```
double *arr = {1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0};
```

Assuming double occupies 8 bytes of memory space, will this line of code work?

```
__m256d first = _mm256_load_pd(arr + 1);
```

- a) True
- b) False

Aligning static arrays

`double c[256];` **May or may not be aligned**

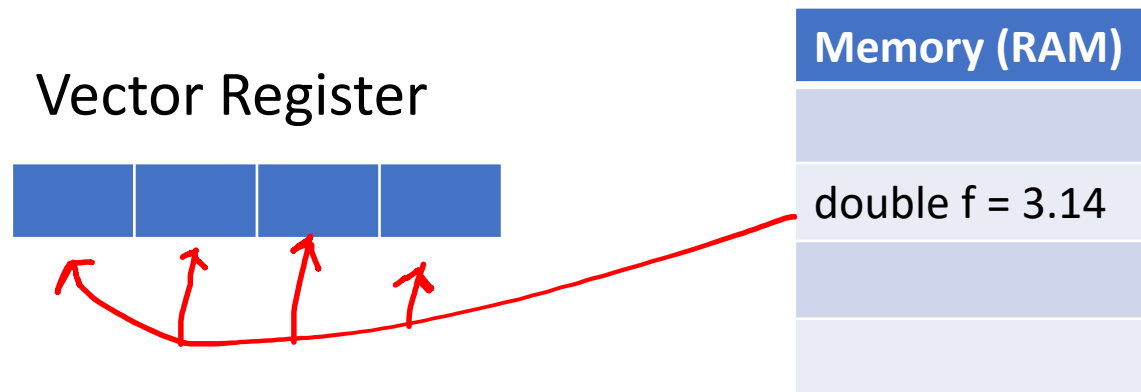
Forcing alignment

```
double c[256] __attribute__((aligned(32)));
```

array c's starting address aligned by 32 bytes

Broadcast

- Loads and broadcast a value to fill a vector register



- How to copy the same value in all the elements of register?
- `_mm256_broadcast_sd(address of double variable)`
returns the vector

```
double *arr = {1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0};  
int i = 2;  
__m256d c0 = _mm256_broadcast_sd( arr + i);
```

Same as:

```
__m256d c0 = _mm256_broadcast_sd( &arr[i]);
```

What does c0 contain?

a) 1.0

b) 1.0, 1.0, 1.0, 1.0

c) 3.0

d) 3.0, 3.0, 3.0, 3.0

Initialize a vector type

```
__m256d c0 = _mm256_set_pd(0.0, 0.0, 0.0, 0.0);
```

Similar to initialization

c0 now contains {0, 0, 0, 0}

c0

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
|---|---|---|---|

Scalar and Vector Instruction

_mm256_broadcast_**sd**

_mm256_load_**pd**

| | | |
|----|-----------------------------------|---------------|
| X: | s stands for Scalar | 1 data |
| | p stands for Packed | multiple data |

Alignment not an issue for scalar instruction

Alignment question 2

Let us assume that arr array is stored in a block of memory with starting location $(3200)_{10}$

```
double *arr = {1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0};
```

Assuming double occupies 8 bytes of memory space, will this line of code work?

```
__m256d first = _mm256_broadcast_sd(arr + 1);
```

- a) True
- b) False

Matrix Multiplication

The general pattern for matrix multiplication is as follows.

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & & \vdots \\ a_{i1} & a_{i2} & a_{i3} & \dots & a_{in} \\ \vdots & \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1j} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2j} & \dots & b_{2p} \\ b_{31} & b_{32} & \dots & b_{3j} & \dots & b_{3p} \\ \vdots & \vdots & & \vdots & & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nj} & \dots & b_{np} \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1j} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2j} & \dots & c_{2p} \\ \vdots & \vdots & & \vdots & & \vdots \\ c_{i1} & c_{i2} & \dots & c_{ij} & \dots & c_{ip} \\ \vdots & \vdots & & \vdots & & \vdots \\ c_{m1} & c_{m2} & \dots & c_{mj} & \dots & c_{mp} \end{bmatrix}$$

$a_{i1}b_{1j} + a_{i2}b_{2j} + a_{i3}b_{3j} + \dots + a_{in}b_{nj} = c_{ij}$

Storing matrix in memory

Assumption: Matrices are stored in column-major order

| | | | |
|----|----|----|----|
| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

Storage is linear – put one column after another column.

1 5 9 13 2 6 10 14 3 7 11 15 4 8 12 16

Matrix Matrix Multiply

- Unoptimized code:

```
1. void dgemm (int n, double* A, double* B, double* C)
2. {
3.   for (int i = 0; i < n; ++i)
4.     for (int j = 0; j < n; ++j)
5.       {
6.         double cij = C[i+j*n];          /* cij = C[i][j] */
7.         for(int k = 0; k < n; k++ )
8.           cij += A[i+k*n] * B[k+j*n];    /* cij += A[i][k]*B[k][j] */
9.         C[i+j*n] = cij;                  /* C[i][j] = cij */
10.      }
11. }
```


X86-64 Registers

- Regular integer registers

Examples

%r11 (64 bits)

%rcx (64 bits)

%rsi (64 bits)

%eax (32 bits)

%ebx (32 bits)

| 63 | 31 | 15 | 7 | 0 | |
|------|-------|-------|-------|---|---------------|
| %rax | %eax | %ax | %al | | Return value |
| %rbx | %ebx | %bx | %bl | | Callee saved |
| %rcx | %ecx | %cx | %cl | | 4th argument |
| %rdx | %edx | %dx | %dl | | 3rd argument |
| %rsi | %esi | %si | %sil | | 2nd argument |
| %rdi | %edi | %di | %dil | | 1st argument |
| %rbp | %ebp | %bp | %bpl | | Callee saved |
| %rsp | %esp | %sp | %spl | | Stack pointer |
| %r8 | %r8d | %r8w | %r8b | | 5th argument |
| %r9 | %r9d | %r9w | %r9b | | 6th argument |
| %r10 | %r10d | %r10w | %r10b | | Caller saved |
| %r11 | %r11d | %r11w | %r11b | | Caller saved |
| %r12 | %r12d | %r12w | %r12b | | Callee saved |
| %r13 | %r13d | %r13w | %r13b | | Callee saved |
| %r14 | %r14d | %r14w | %r14b | | Callee saved |
| %r15 | %r15d | %r15w | %r15b | | Callee saved |

Figure 3.2 Integer registers. The low-order portions of all 16 registers can be accessed as byte, word (16-bit), double word (32-bit), and quad word (64-bit) quantities.

How to read a vector instruction?

Format nameXY

X: **s** stands for **Scalar**

p stands for **Packed**

Y: Y is for data type. We will use only double so Y = d

vaddsd : **scalar** double

vaddpd : **packed** double

vmulsd : **scalar** double

vmulpd : **packed** double

Matrix Multiply

- x86 assembly code:

```
1. vmovsd (%r10),%xmm0    # Load 1 element of C into %xmm0
2. mov %rsi,%rcx          # register %rcx = %rsi
3. xor %eax,%eax          # register %eax = 0
4. vmovsd (%rcx),%xmm1    # Load 1 element of B into %xmm1
5. add %r9,%rcx           # register %rcx = %rcx + %r9
6. vmulsd (%r8,%rax,8),%xmm1,%xmm1    # Multiply %xmm1, element of A
7. add $0x1,%rax          # register %rax = %rax + 1
8. cmp %eax,%edi          # compare %eax to %edi
9. vaddsd %xmm1,%xmm0,%xmm0    # Add %xmm1, %xmm0
10. jg 30 <dgemm+0x30>    # jump if %eax > %edi
11. add $0x1,%r11d        # register %r11 = %r11 + 1
12. vmovsd %xmm0, (%r10)  # Store %xmm0 into C element
```

Matrix Multiply

- Optimized C code:

```
1. #include <immintrin.h>
2. void dgemm (int n, double* A, double* B, double* C)
3. {
4.     for ( int i = 0; i < n; i+=4 )
5.         for ( int j = 0; j < n; j++ ) {
6.             __m256d c0 = _mm256_load_pd(C+i+j*n); /* c0 = C[i][j] */
7.             for( int k = 0; k < n; k++ )
8.                 c0 = _mm256_add_pd(c0, /* c0 += A[i][k]*B[k][j] */
9.                                     _mm256_mul_pd(_mm256_load_pd(A+i+k*n),
10.                                                    _mm256_broadcast_sd(B+k+j*n)));
11.             _mm256_store_pd(C+i+j*n, c0); /* C[i][j] = c0 */
12.         }
13. }
```

Matrix Multiply

- Optimized x86 assembly code:

```
1. vmovapd (%r11),%ymm0      # Load 4 elements of C into %ymm0
2. mov %rbx,%rcx             # register %rcx = %rbx
3. xor %eax,%eax             # register %eax = 0
4. vbroadcastsd (%rax,%r8,1),%ymm1 # Make 4 copies of B element
5. add $0x8,%rax             # register %rax = %rax + 8
6. vmulpd (%rcx),%ymm1,%ymm1 # Parallel mul %ymm1,4 A elements
7. add %r9,%rcx              # register %rcx = %rcx + %r9
8. cmp %r10,%rax             # compare %r10 to %rax
9. vaddpd %ymm1,%ymm0,%ymm0  # Parallel add %ymm1, %ymm0
10. jne 50 <dgemm+0x50>      # jump if not %r10 != %rax
11. add $0x1,%esi            # register % esi = % esi + 1
12. vmovapd %ymm0, (%r11)    # Store %ymm0 into 4 C elements
```

Comparison

| Unoptimized | Optimized |
|--|--|
| <p>Scalar double (sd) instructions</p> <pre>vmovsd (%r10), %xmm0</pre> <p># Load 1 element of C into %xmm0</p> <pre>vmulsd</pre> <p>xmm register = 128 bits wide</p> | <p>Parallel double (pd) instructions</p> <pre>vmovapd (%r11), %ymm0</pre> <p># Load 4 elements of C into %ymm0</p> <pre>vmulpd</pre> <p>ymm register = 256 bits wide</p> <p>3 to 4 times faster by using Intrinsics</p> |

HW: GFLOPS calculation

- Use matmul_intrinsics.c for testing your code
 - It has the logic and code for timing and gflops.
 - testing is same as matrix multiplication, you can reuse the code for array creation, initialization.

```
double time_spent = (double) (tv2.tv_usec - tv1.tv_usec) / 1000000 +  
                    (double) (tv2.tv_sec - tv1.tv_sec);  
  
printf("Intrinsics: Time %f \n", time_spent);  
  
double numOps = 2.0 * SIZE * SIZE * SIZE;  
  
double gflops = 1.0e-9 * numOps / time_spent;  
  
printf(" Intrinsics: numOps = %.2f GFLOPS %f \n", numOps, gflops);
```


- Code-reuse from matmul intrinsics code

```
int matmul_intrinsics()
{
    // int NUM_ELEMENTS = 1024*1024;
    // int dimension = 1024;

    int NUM_ELEMENTS = SIZE*SIZE;
    int dimension = SIZE;
    printf("dimension = %d & NUM_ELEMENTS %d, sizeof(double)=%lu \n ", dimension, NUM_ELEMENTS, sizeof(double));

    size_t N_pd = (NUM_ELEMENTS*8)/sizeof(double);

    double *data_A = (double*)_mm_malloc(N_pd*sizeof(double), 32);
    double *data_B = (double*)_mm_malloc(N_pd*sizeof(double), 32);
    double *data_C = (double*)_mm_malloc(N_pd*sizeof(double), 32);

    if(data_A == NULL || data_B == NULL || data_C == NULL)
    {
        printf("Error \n");
        return 1;
    }

    initialize(data_A, NUM_ELEMENTS, 1.0);
    //initialize(data_B, NUM_ELEMENTS, 2.0);
    initializeMM(data_B, NUM_ELEMENTS);
    initialize(data_C, NUM_ELEMENTS, 0.0);

    struct timeval tv1, tv2;

    gettimeofday(&tv1, NULL);

    //dgemmIntrin(dimension, data_A, data_B, data_C);

    vecAddIntrin(dimension, data_A, data_B, data_C);

    gettimeofday(&tv2, NULL);

    double time_spent = (double) (tv2.tv_usec - tv1.tv_usec) / 1000000 +
        (double) (tv2.tv_sec - tv1.tv_sec);

    printf("Intrinsics: Time %f \n", time_spent);
}
```

Compiling on Pascal

- Use Pascal server (Intel Processor required)
- Header file required
`#include<immintrin.h>`
- `gcc -O3 -mavx arrayIntrinsics.c`
OR
`gcc -O3 -o prog -mavx matmul_intrinsic.c`
- `./a.out`

Vectorization by compilers

- Compiler options/flags

`gcc -O2`

`gcc -O3`

- Vectorization is parallelization.
- Look at assembly for vector instructions
- Compiler vectorization report
- Compiler may do sub-optimal vectorization