# TDT4200 Parallel Computing
# Problem Set 3

Stian Hvatum (hvatum)
MTDT

September 23, 2012

## Contents

## Task 1 Theory

### a) Reduce branches

1. If we are using if's with elsif's, we can put the most common case first, so that we usually only do one branch. Code Listing a) 1.

2. Some conditional statements can be sustituted with equals that are not conditional. Code listing a) 2.

3. We can do loop-unrolling. If we know number of repetitions is known to be a multiple of N>1, this can easily be done. Code listing a) 3.

### a) 1.

```
1   for (int i = 0; i < N; i++)
2   {
3     if (i == N/2 + 1)
4     {
5       do_seldome_stuff();
6     }
7     else if (i % N/4 == 1)
8     {
9       do_stuff();
10    }
11    else if (i % 2 == 0)
12    {
13      do_common_stuff();
14    }
15  }
```

Will be faster if we put the common stuff first (less checking). This implies that the conditions are exclusive.

```
1   for (int i = 0; i < N; i++)
2   {
3     if(i % 2 == 0)
4     {
5       do_common_stuff();
6     }
7     else if(i == N/2 + 1)
8     {
9       do_seldome_stuff();
10    }
11  // This condition will be more common,
12  // but have overlaping conditions with previous statement
13    else if(i % N/4 == 1)
14    {
15      do_stuff();
16    }
17  }
```

### a) 2.

We can substitute this static conditional assignment with some clever logic statement.

```
1    if (n % 2 == 0)
2        a[n] = 2;
3    else
4        a[n] = 0;
```

```
1        a[n] = (n & 0b) << 1;
```

### a) 3.

In this example, N/2 is allways a multiple of 2 (since we are dividing in half, which gives an even number).

```
1    for (int i = 0; i < N/2; i++)
2    {
3        do_stuff(i);
4    }
```

Since i is known to be a multiple of 2, we can unroll the loop at least twice, like shown here.

```
1    for (int i = 0; i < N/2; i++)
2    {
3        do_stuff(i);
4        do_stuff(++i);
5    }
```

### b) Improve switches

We can improve the performance of a long switch by placing the most common case first. We can also splitt it into nested switch statements or nest inside the else-block of an if-else-statement where the true-block of the if-statement contains the most frequent cases.

**Example:** Let us assume *answer()*'s most common return value is 42, but can return a lot of other values, then this statement will have improved performance over a regular switch statement with sequentially ordered cases. Of cource, if this switch is converted into a jump-table, or 42 was the first case AND the switch was converted into if-else-statements by the compiler, this optimization won't matter that much, as it would be implicit.

```
1    int a = answer();
2    if (a == 42)
3        hurray();
4    else
5        switch (a)
6        {
7        case 1:
8            ...
9            break;
10       case 2:
11           ...
12           break;
13       ...
14       ...
15       default:
16           blah();
17           break;
18       }
```

### c) Reduce cache misses

1. One can change loops to operate on blocks instead of rows of data. This may help to reduce directe-mapped cache conflicts.

2. Two arrays that are concurrently accessed can be merged with an array of structs, where the struct contains one element from each original array. This helps if both elements are needed, and the original arrays conflicts in cache.

3. If we are looping though a 2-dim array, and we have the first index looped over at the inner loop, we can switch the order of the loops, such that the inner loop controlls the second index. This helps since the first index effects a stride of the size of the inner array (or random if malloc'ed in C), while the second index gives a stride of 1 element.

### d) Implementing SIMD

#### 1. As loop unrolling

In its simplest form, vectorization can be used as a form of loop unrolling, as we do two and two double value arithmetics at the same time. This will speed up the loop, and is fairly simple. Disadvantages is that vectorization in general may increase overhead in flow control of the program[1].

#### 2. For complex artithmetics

Intel has a lot of docs on how to do complex artithmetics with SSSE3. Eg. can you do complex multiplication in just a few steps, while the C99 Complex library implements a rather large routine for this purpose. SSSE3 also has an *addsub* instruction which is ideal for the last step of complex multiplication, where we add and subtract in the same operation. This is theoretically much faster, but code readability are greatly reduced.

#### 3. General math

When we are facing math where multiple steps involve the same values, we can use SIMD to speed up the process. This might be tricky and hard to read, but might also speed up the program.

#### 4. Let the compiler do it

When specifyng -O3 -march=native, and your CPU supports SSSE3, GCC will use it where it finds it usefull.

## Task 2 Code optimization

### Task 2.1 Optimizing Taylor series for Sinus approximation

#### 1. Less branches

One obvious optimization we can do[2], is that the *factorial* could be done iterative with a loop, instead of recursive calls. This saves a lot of function calls, which contains a lot

---

[1] http://en.wikipedia.org/wiki/Vectorization_(parallel_computing)#Reducing_vectorization_overhead_in_the_presence_of_control_flow

[2] At first look, I thought the function *factorial*did not need to return a double, but since its values are too large, integers would overflow, thus we need the double value.

of stack- and register-maintenance. The *slow_sin* is doing 100 branches with different outcome at each repetition, this is bad for branch prediction. We can modify the code to instead of checking and doing two different operations, we can do the same thing each time. On my computer, these two optimizations gives the function some more than double speed[3].

```
1   /* calculate n! */
2   double factorial(int n)
3   {
4       double answer = n;
5       for (int i = n-1; i > 1; i--)
6           answer *= i;
7       return answer;
8   }
9   /* calculate sin(x) using Taylor series expansion */
10  double slow_sin(double x)
11  {
12      double r=0; int i;
13      for(i=0;i<100;i++)
14      {
15      int sign = (i & 1) * -2 + 1;
16      r += pow(x, i*2+1) / factorial(i*2+1) * sign;
17      }
18  return r;
19  }
```

## 2. Loop unrolling

Instead of doing the sign test for + or -, we could simply unroll the loop, being shure each even step is + and each odd step is -. We could unroll the loop as many times we want, effect is determined by cache size, since we want to keep the inner instructions in cache, instead of fetching them from a long distance for each entry. With this optimization, along with the optimized *factorial*, speed is about 3x of the original[4].

```
1
2   /* calculate sin(x) using Taylor series expansion */
3   double slow_sin(double x)
4   {
5       double r=0; int i;
6       for(i=0;i<100;i++)
7       {
8       r += pow(x, i*2+1) / factorial(i*2+1);
9       i++;
10      r -= pow(x, i*2+1) / factorial(i*2+1);
11      }
12  return r;
13  }
```

## 3. Using threads

Since the result is not dependent of the order of the partial results, we can use pthreads, MPI or similar to split the problem into smaller problems. In practice, this slowed things down at my computer.

---

[3]tested with *time* and setting the sine-loop to 10000 instead of 100. Results was new function using 0.4 sec, old function using 0.9 sec. (Intel Core i7 930@ 2.8GHz, GCC -O0)

[4]0.32 sec. on my machine, same configuration as previous test

```
1    ...
2    /* calculate sin(x) using Taylor series expansion */
3    double slow_sin_new(double x, int start, int stride)
4    {
5        double r=0; int i;
6        for(i=start;i<10000;i+stride)
7        {
8        r += pow(x, i*2+1) / factorial_new(i*2+1);
9        i++;
10       r -= pow(x, i*2+1) / factorial_new(i*2+1);
11       }
12   return r;
13   }
14   ...
15   void* start_thread(void* rank)
16   {
17       double* ans = malloc(sizeof(double));
18       *ans = slow_sin_new(input, (int)rank * 2, 15);
19       pthread_exit((void*)ans);
20   }
21   ...
22   int main()
23   {
24       pthread_t threads[8];
25       for (int i = 0; i < 8; i++)
26           pthread_create(&threads[i], NULL, start_thread, (void*)i);
27       for (int i = 0; i < 8; i++)
28       {
29           double val;
30           pthread_join(threads[i], (void*)&val);
31           r += val;
32       }
33       printf("%f\n", r);
34   }
```

## 4. Using SSSE3

So instead of threads, we can use vectorization via SSSE3 or equal instruction set. This function uses pow and factorial, but we can also implement them inline with SSSE3 instructions. Speedup was disapointingly only double from original[5], but this is most likely since I am calling pow and factorial and stores the result into special registers inbetween the other SSSE3 instructions. This can most likely be avoided by inlining these two functions with SSSE3 instructions[6].

```
1    /* calculate sin(x) using Taylor series expansion */
2    double slow_sin_new(double x)
3    {
4        int i; double r = 0;
5        __m128d factor, res, part, cx, nx;
6        double xarr[2];
7        res = _mm_set_pd(0,0);
8
9        for(i=0;i<10000;i+=2)
10       {
11       factor = _mm_set_pd(factorial(i*2+1), factorial(i+2+3));
12       part = _mm_set_pd(pow(x, i*2+1), pow(x, i+2+3));
13       res = _mm_addsub_pd(res, part);
14       }
15       _mm_storeu_pd(xarr, res);
16       return xarr[0] + xarr[1];
17   return r;
18   }
```

---

[5]0.45 seconds

[6]Pow can be sustituted with a multiplication with x for each step, and factorial can be a multiplication with i at each step.

**Task 2.2 Loops and indexed access optimizations**

As explained in Task 1 c) 3., this code access the a-array with the first index in the inner loop, which gives a large stride in the accessed data. This gives the cache a hard time keeping up, since accessed data is not aligned in memory the same way it is accessed. Set loop over $j$ as outer loop, and over $i$ as inner loop in order to speed up function. This loop could also benefit from being unrolled and/or vectorized.