# Parallel Programming SS21 Final Project

Project 02 – Gaussian Elimination
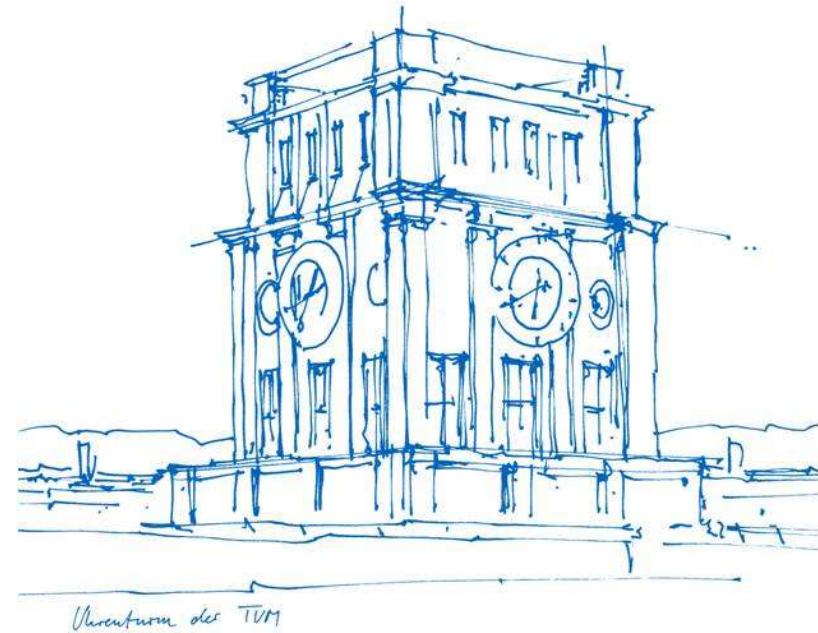
Group 215

09.07.2021

Zhelin Yang

Dian Yuan

Jing Xiong

TUM

Uhrenturm der TUM

# Sequential code analysis  - Profiling

### size 1024x1024



### size 2048x2048



### size 4096x4096



### size 8192x8192

# Sequential code analysis  - Profiling

```
Performance counter stats for './serialge ./ge data/size8192x8192':

     216.171,10 msec task-clock                #    1,000 CPUs utilized
          3.996      context-switches          #    0,018 K/sec
              2      cpu-migrations            #    0,000 K/sec
        131.257      page-faults               #    0,607 K/sec
893.168.658.549      cycles                    #    4,132 GHz
2.459.203.457.971    instructions              #    2,75  insn per cycle
565.625.107.143      branches                  # 2616,562 M/sec
     42.950.841      branch-misses             #    0,01% of all branches

   216,187621527 seconds time elapsed

   216,013328000 seconds user
     0,160030000 seconds sys
```

# Sequential code analysis  -  Amdahl's law

TUM

Parameters:

f = fraction of parallel execution                f = 99.8%(size2048x2048)

p = number of parallel tasks/threads/processes        p = 4

Speedup refers to speedup of function "Solve"

Theoretical Speedup:

$$SU(p) = \frac{T}{T(p)} = \frac{T}{(1-f)*T + \frac{f*T}{p}} = \frac{1}{1-f+\frac{f}{p}} = \frac{1}{1-0.998+\frac{0,998}{4}} = 3.976$$

# Sequential code analysis  -  Amdahl's law

Parameters:

f = fraction of parallel execution                           f ≈ 100%

p = number of parallel tasks/threads/processes               p = 4

Speedup refers to speedup of function "Solve"

$$SU(p) = \frac{T}{T(p)} = \frac{T}{(1-f)*T + \frac{f*T}{p}} = \frac{1}{1 - f + \frac{f}{p}} = \frac{1}{1 - 1 + \frac{1}{4}} = 4$$

**Maximal Speedup: Around 4**

# OpenMP - Parallelized implementation and approach

```c
void ForwardElimination(double *matrix, double *rhs, int rows, int columns){
    for(int row = 0; row < rows; row++){
        // Extract Diagonal element
        int diag_idx = row*rows + row;
        double diag_elem = matrix[diag_idx];
        #pragma omp parallel for schedule(dynamic)
        for (int lower_rows=row+1; lower_rows<rows; lower_rows++){
            int below_diag_idx = lower_rows*rows + row;
            assert(diag_elem!=0);
            // Compute the factor
            double elimination_factor = matrix[below_diag_idx]/diag_elem;
            int element_idx;
            for (int column=row+1; column<columns; column++){
                // set the column index of the entry to be operated
                element_idx = lower_rows*rows + column;
                // subtract the row
                matrix[element_idx] -= elimination_factor*matrix[row*rows+column];
            }
            rhs[lower_rows] -= elimination_factor*rhs[row];
            // set below diagonal elements to 0
            // matrix[below_diag_idx] = 0.;
        }
    }
}
```

**First try**

-Totally 3 „for loops" in this function

-Parallel region: 2nd „for loop"

-Schedule:dynamic

# OpenMP - Intermediate Speed-up results, profiling

| size | Overhead of ForwardElimination | Overhead of vmovapd | Speedup (on my PC) | Speedup (on server) |
|------|-------------------------------|---------------------|--------------------|---------------------|
| 2048 | 80.99% | 69.15% | 1.39 | ≈10 |
| 4096 | 90.44% | 73.9% | 1.23 | |
| 8192 | 95.06% | 75.82% | 1.14 | |

PS: speedup on PC is a speedup for function "Solve" compared with the optimized sequential code

**Result of first try**

-Bottleneck: ForwardElimination--Data Transfer

-Speedup: Not bad on my PC, but on the server it is not fast enough

# OpenMP - Parallelized implementation and approach

□ □ □ □ □ □ □ □ □  ← row

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥  thread 0

🟨 🟨 🟨 🟨 🟨 🟨 🟨 🟨  thread 1

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩  thread 2

🟨 🟨 🟨 🟨 🟨 🟨 🟨 🟨  thread 1

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥  thread 0

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩  thread 2

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥  thread 0

**First try**

-Distribution: The task row addresses distributed to one thread are not continuous

# OpenMP - Parallelized implementation and approach



← row

thread 2
thread 1
thread 0
thread 1
thread 0
thread 2

**First try**

-Distribution: The task row addresses distributed to one thread are not continuous and randomly

# OpenMP - Final Implementation improvements and new speed-up

```cpp
void ForwardElimination(double *matrix, double *rhs, int rows, int columns){
    for(int row = 0; row < rows; row++){
        // Extract Diagonal element
        int diag_idx = row*rows + row;
        double diag_elem = matrix[diag_idx];
        #pragma omp parallel for
        for (int lower_rows=row+1; lower_rows<rows; lower_rows++){
            int below_diag_idx = lower_rows*rows + row;
            assert(diag_elem!=0);
            // Compute the factor
            double elimination_factor = matrix[below_diag_idx]/diag_elem;
            int element_idx;
            for (int column=row+1; column<columns; column++){
                // set the column index of the entry to be operated
                element_idx = lower_rows*rows + column;
                // subtract the row
                matrix[element_idx] -= elimination_factor*matrix[row*rows+column];
            }
            rhs[lower_rows] -= elimination_factor*rhs[row];
            // set below diagonal elements to 0
            // matrix[below_diag_idx] = 0.;
        }
    }
}
```

**Last Version**

-Schedule:default

-More server friendly

# OpenMP - Parallelized implementation and approach

□ □ □ □ □ □ □ □ □ ← row

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥 thread 0

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥 thread 0

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥 thread 0

🟨 🟨 🟨 🟨 🟨 🟨 🟨 🟨 thread 1

🟨 🟨 🟨 🟨 🟨 🟨 🟨 🟨 thread 1

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩 thread 2

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩 thread 2

**Last version**

-The thread will access similar address in the next round

# OpenMP - Parallelized implementation and approach

TUM

row ←

thread 0

thread 0

thread 1

thread 1

thread 2

thread 2

## Last version

-The thread will access similar address in the next round

# OpenMP - Final Speed-up results, profiling

| size | Cache misses | |
|---|---|---|
| | dynamic | default |
| 2048 | 57% | 57% |
| 4096 | 84% | 83% |
| 8192 | 90% | 86% |

# OpenMP - Final Speed-up results

| size | Overhead of ForwardElimination | Overhead of vmovapd | Speedup (on my PC) | Speedup (on server) |
|------|-------------------------------|---------------------|--------------------|--------------------|
| 2048 | 80.99% | 69.15% | 1.39 | ≈10(dynamic) ≈14(default) |
| 4096 | 90.44% | 73.9% | 1.23 | |
| 8192 | 95.06% | 75.82% | 1.14 | |

## Result of Final

-Speedup: Similar as the previous version on my PC, but better on server

# OpenMP - Final Implementation

Theoretical speedup: 3~4(with 4 threads) according to Amdahl's law
- Almost all of ForwardElimination is parallized

Bottleneck: Large cost from data transfer
- The bigger the matrix is, the larger the cost from data transfer would be.

# MPI - Parallelized implementation and approach

rank 0
rank 0
rank 1
rank 1
rank 2
rank 2
rank 3
rank 3

**First try**

-Data are equally distributed to each progress in block

-Addresses of data are continuous in each progress

# MPI - Parallelized implementation and approach

row ▢ ▢ ▢ ▢ ▢ ▢ ▢ ▢ ▢    rank 0

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥    rank 0

🟧 🟧 🟧 🟧 🟧 🟧 🟧 🟧    rank 1

🟧 🟧 🟧 🟧 🟧 🟧 🟧 🟧    rank 1

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩    rank 2

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩    rank 2

🟦 🟦 🟦 🟦 🟦 🟦 🟦 🟦    rank 3

🟦 🟦 🟦 🟦 🟦 🟦 🟦 🟦    rank 3

**First try**

-Data are equally distributed to each progress in block

-Addresses of data are continuous in each progress

# MPI - Parallelized implementation and approach

row

rank 0
rank 0
rank 1
rank 1
rank 2
rank 2
rank 3
rank 3

**First try**

-Data are equally distributed to each progress in block

-Addresses of data are continuous in each progress

-Not balanced

# MPI - Final Implementation improvements and new speed-up

rank 0

rank 1

rank 2

rank 3

rank 0

rank 1

rank 2

rank 3

**Final version**

-Data are equally distributed to each progress

-Addresses of data are not continuous in each progress

# MPI - Final Implementation improvements and new speed-up

TLIT

row ▢ ▢ ▢ ▢ ▢ ▢ ▢ ▢ ▢   rank 0

🟨 🟨 🟨 🟨 🟨 🟨 🟨 🟨   rank 1

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩   rank 2

🟦 🟦 🟦 🟦 🟦 🟦 🟦 🟦   rank 3

🟥 🟥 🟥 🟥 🟥 🟥 🟥 🟥   rank 0

🟨 🟨 🟨 🟨 🟨 🟨 🟨 🟨   rank 1

🟩 🟩 🟩 🟩 🟩 🟩 🟩 🟩   rank 2

🟦 🟦 🟦 🟦 🟦 🟦 🟦 🟦   rank 3

**Final version**

-Data are equally distributed to each progress

-Addresses of data are not continuous in each progress

# MPI - Final Implementation improvements and new speed-up

row

| | rank 0 |
| | rank 1 |
| | rank 2 |
| | rank 3 |
| | rank 0 |
| | rank 1 |
| | rank 2 |
| | rank 3 |

**Final version**

-Data are equally distributed to each progress

-Addresses of data are not continuous in each progress

-Almost balanced

# MPI – Speedup

| size | Speedup on my PC | | Speedup Server | |
|------|------------------|----------------|----------------|----------------|
|      | First try        | Final version  | First try      | Final version  |
| 2048 | 1.39             | 1.39           | 10.54          | 10.89          |
| 4096 | 1.21             | 1.18           |                |                |
| 8192 | 1.13             | 1.12           |                |                |

**Comparison**

-Speedup: Very close both on PC and server

-The impact of imbalance is not so obvious

# MPI - Parallelized implementation and approach

```
# Total Lost Samples: 0
#
# Samples: 2M of event 'cycles'
# Event count (approx.): 1016468430965
#
# Overhead        Pid:Command
# ........    ....................
#
    25.05%      7329:mpige
    25.04%      7330:mpige
    25.03%      7332:mpige
    24.88%      7331:mpige
     0.00%      7328:hydra_pmi_proxy
     0.00%      7326:mpirun
```

## Comparison

- The impact of imbalance is not so obvious

- The load seems balanced because there is a barrier at each end of the first "for loop" to update "row+1"

```
MPI_Bcast(matrix+(row+1)*rows, columns, MPI_DOUBLE, r, MPI_COMM_WORLD);
MPI_Bcast(rhs+row+1, 1, MPI_DOUBLE, r, MPI_COMM_WORLD);
```

# MPI - Profiling

| size | First try | | Final version | |
|------|-----------|---|---------------|---|
| | ForwardElimination | vmovapd | ForwardElimination | vmovapd |
| 2048 | 45.13% | 52.3% | 53.37% | 67.27% |
| 4096 | 57.31% | 56.98% | 73.14% | 72.68% |
| 8192 | 64.01% | 57.75% | 85.03% | 74.77% |

## Comparison

- The data address in the first try are continuous, so that it would be easier to access

# MPI - Final Speed-up results, profiling

TITI



| 64.01% | mpige | mpige | [.] | MPI::Forw |
|--------|-------|-------|-----|-----------|
| 3.80% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 1.12% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.99% | mpige | [kernel.kallsyms] | [k] | do_syscal |
| 0.95% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.82% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.81% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.81% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.81% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.80% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.80% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |
| 0.80% | mpige | libmpich.so.0.0.0 | [.] | 0x0000000 |

**Cost from MPI**

-Using MPI might cause some communication cost

# MPI - Final Implementation

Theoretical speedup: 3~4(with 4 threads) according to Amdahl's law
- Almost all of ForwardElimination is parallized

Bottleneck: Large cost from data transfer, MPI communication
- The bigger the matrix is, the larger the cost from data transfer would be
- MPI communication would cost some time

# Hybrid - Parallelized implementation and approach

Combine of OMP and MPI approach

- Two threads for MPI(distributed nodes) and two for OMP(multiple cores)
- The structure is almost the same as MPI approach
- OMP is nested within the MPI process

```
void ForwardElimination(double *matrix, double *rhs, int rows, int columns, int rank, int size){
    for(int row = 0; row < rows-1; row++){
        int block_idx = row / size;
        int l_rows_start;

        if((row%size)<rank){
            l_rows_start = block_idx*size + rank;
        }

        else{
            l_rows_start = (block_idx+1)*size + rank;
        }

        ...

        #pragma omp parallel for num_threads(2)
        for (){
            ...
        }

        ...

        int r = (row+1)%size;
        MPI_Bcast(matrix+(row+1)*rows, columns, MPI_DOUBLE, r, MPI_COMM_WORLD);
        MPI_Bcast(rhs+row+1, 1, MPI_DOUBLE, r, MPI_COMM_WORLD);
    }
}
```

Master thread

Worker threads of the master processor

End OpenMP

# Hybrid - Final Performance Results

| size | Speedup |
|------|---------|
| 2048 | 1.39 |
| 4096 | 1.20 |
| 8192 | 1.13 |

*with 4 threads on own PC*

```
#  O v e r h e a d          P i d : C o m m a n d
#  . . . . . . . .          . . . . . . . . . . . . . . . . . . . . . . .
#
      2 7 . 6 1 %              4 1 2 2 : h y b r i d g e
      2 6 . 9 9 %              4 1 2 3 : h y b r i d g e
      1 8 . 2 0 %              4 1 2 5 : h y b r i d g e
      1 8 . 0 0 %              4 1 2 4 : h y b r i d g e
```

The two threads responsible for the MPI process share more load

Theoretical speedup: 3~4 (with 4 threads) according to Amdahl's law
- Almost all of ForwardElimination is parallized

Bottleneck: Large cost from data transfer, MPI communication
- The bigger the matrix is, the larger the cost from data transfer would be
- MPI communication would cost some time
- Load balance

**-But this is a simulation on a single machine, it should perform better on real distribution nodes**

# Bonus (SIMD) - Parallelized implementation and approach ◻◻◻

```
for (int column=row+1; column<columns; column++){
    // Set the column index of the entry to be operated
    element_idx = lower_rows*rows + column;
    // Subtract the row
    matrix[element_idx] -= elimination_factor*matrix[row*rows+column];
}
```

This is operation on the row and can be vectorized.

AVX Intrinsics



_m256d (four 64-bit float)

The rest do the normal operation

```
int vectorLen = 4;
int end = columns - columns % vectorLen;

__m256d elimination_factor_vec=_mm256_set1_pd(elimination_factor);
__m256d element_row, select_row;

for(int column = row + 1;column < end;column += vectorLen){
    element_row = _mm256_loadu_pd(matrix + lower_rows*rows + column);
    select_row = _mm256_loadu_pd(matrix + row*rows + column);
    __m256d temp_row = _mm256_mul_pd(elimination_factor_vec, select_row);
    element_row = _mm256_sub_pd(element_row, temp_row);
    _mm256_storeu_pd(matrix + lower_rows*rows + column, element_row);
}

//Eliminate the rest elements in the row
for(int column = end; column < columns; column++){
    matrix[lower_rows*rows + column] -= elimination_factor*matrix[row*rows+column];
}
rhs[lower_rows] -= elimination_factor*rhs[row];
}
```

# Bonus (SIMD) – Final Performance Results

| size | Hybrid Speedup | Bonus Speedup *with 4 threads on own PC* |
|------|----------------|------------------------------|
| 2048 | 1.39 | 1.40 |
| 4096 | 1.20 | 1.20 |
| 8192 | 1.13 | 1.13 |

No significant improvement on speedup.

**Let's look at assembly of *hybridge.cpp***

```
.L8:
    vmovupd (%rax,%r12), %xmm6
    vinsertf128 $0x1, 16(%rax,%r12), %ymm6, %ymm1
    vmovupd 0(%r13,%r12), %xmm5
    vinsertf128 $0x1, 16(%r13,%r12), %ymm5, %ymm0
    vmulpd  %ymm3, %ymm1, %ymm1
    vsubpd  %ymm1, %ymm0, %ymm0
    vmovups %xmm0, 0(%r13,%r12)
    vextractf128    $0x1, %ymm0, 16(%r13,%r12)
    addq    $32, %r12
    cmpq    %rsi, %r12
    jne .L8
```

Many vectorized AVX/SSE instructions have been generated after compilation.
It's the characteristic of gcc -o3

# Conclusion

*Speedup with 4 threads on own PC*

| size | OpenMP | MPI | OMP+MPI | OMP+MPI+SIMD |
|------|--------|------|---------|--------------|
| 2048 | 1.39   | 1.39 | 1.39    | 1.40         |
| 4096 | 1.23   | 1.18 | 1.20    | 1.20         |
| 8192 | 1.14   | 1.12 | 1.13    | 1.13         |

- Speedup are close among the four approaches.
- Speedup doesn't reach the theoretical speedup mostly due to cost of data transfer
- It would be more appropriate to use MPI and hybrid method for distributed system