# CHALMERS
## EXAMINATION / TENTAMEN

| Course code/kurskod | Course name/kursnamn | | | |
|---|---|---|---|---|
| DIT431 | High Performance Parallel Programming | | | |
| Anonymous code Anonym kod | | Examination date Tentamensdatum | Number of pages Antal blad | Grade Betyg |
| DIT431-0007-ADT | | 27/10/2023 | 13 | |

I confirm that I've no mobile or other similar electronic equipment available during the examination.
Jag intygar att jag inte har mobiltelefon eller annan liknande elektronisk utrustning tillgänglig under eximinationen.

| Solved task Behandlade uppgifter No/nr | | Points per task Poäng på uppgiften | Observe: Areas with bold contour are to completed by the teacher. Anmärkning: Rutor inom bred kontur ifylles av lärare. |
|---|---|---|---|
| 1 | X | 8 | |
| 2 | X | 8 | |
| 3 | X | 2 | |
| 4 | X | 9 | |
| 5 | X | 9 | |
| 6 | X | 6 | |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |
| 11 | | | |
| 12 | | | |
| 13 | | | |
| 14 | | | |
| 15 | | | |
| 16 | | | |
| 17 | | | |
| Bonus poäng | | | |
| Total examination points Summa poäng på tentamen | | 42 | |

**CHALMERS**

Anonymous code / Anonym kod

DIT431-0007-ADT

Points for question (to be filled in by teacher)
Poäng på uppgiften (ifylles av lärare)

Consecutive page no. / Löpande sid nr **1**

Question no. / Uppgift nr **1**

(a) ~~(i) N:1.~~
~~1000 ULTs.~~

~~$\frac{1000}{10} = 100$ ULTs on each core.~~

~~Total execution time = $100 * (10 + 1)$~~
~~= 1100 units~~

(i) N:1:
(a) Only 1 core can be used.

1000 ULTs.
Total execution time = $1000 * (10+1)$
$= 11000$ units.

(2) 1:1

$\frac{1000}{10} = 100$ ~~ULTs~~ KLTS on each core.

Total execution time = $100(10 + 5)$
$= 1500$ units.

(3) N:M, 10:1

No. of KLTs = $\frac{1000}{10} = 100$.

~~total~~ execution time = $(10) *$ time for 1 ULT
of 1 KLT
$= 10 * (10 + 1)$
$= 110$ units

total execution time = $\frac{100}{10} (110 + 5)$
$= 1150$ units.

1b) Each task now takes 20 units, thus the execution time to context switch time ratio increases. to 20:1 for ULTs and 20:5 for KLTs.

This makes using ~~KLT~~ ~~MLT more suitable~~ more fine grained threading models more suitable.

~~N:1 results in execution time~~

N:1 is not suitable as no parallelism is utilised.

1:1 has execution time of $100(20+5) = 2500$ units.

10:1 has $10 \times (5 + 10(1+20)) = 2150$ units.

100:1 has $1 \times (100(1+20)) = 2100$ units

Thus 100:1 is most effective.

**2.5**

1c) If there are more than 1 KLT like in 1:1 or N:M, the KLTs have a shared address space.

**0.5**
**0**
Thus having multiple KLTs should have no issues on the program, as each KLT is able to pass the message to the different node.

(d) CUDA threads are significantly more fine grained. ~~Cuda~~ CUDA threads have private registers but shared memory between threads in the same block. There is also global memory that all threads can ~~can~~ access. There are many CUDA cores that can execute the many fine grained CUDA threads. Context switching overheads are also low as each thread has its own registers. Similar to SIMD, CUDA uses Single Instruction Multiple Threads by executing ~~finding~~ the same instruction on all threads in a warp.

**2.5**

**CHALMERS**

5

| Anonymous code | Points for question | Consecutive page no. |
| | (to be filled in by teacher) | Löpande sid nr |
| Anonym kod | Poäng på uppgiften | Question no.    3 |
| | (ifylles av lärare) | Uppgift nr    2 |

DIT431-0007-ADT

2 a)   a) Static scheduling.
64 loop iterations. Each processor has $\frac{64}{8} = 8$ iterations.

Time units to complete = 8 × 2.0
= 16 time units.    ×64

2a (b)   Self-scheduling.

Each processor has 8 iterations

time units to complete = 8 × 2.0 = 16 time units.    ×64

2a (c)   Chunk-scheduling.

$\frac{64}{6}$ = ~~16 chunks~~. 10.6 chunks

~~Each processor executes $\frac{16}{8}$ = 2 chunks.~~

There are 10 chunks of size 6, 1 chunk of size 4.

6      6      4.
6      6      6    6    6    6    6    6
$P_1$   $P_2$   $P_3$  $P_4$  $P_5$  $P_6$  $P_7$  $P_8$

Execution time = (2×6) × 2.0 = 24 time units.    ×64

2a)(d)   guided-scheduling.

(3)

4      4      4      4      4      2
$\frac{64}{8}=8$  $\frac{56}{8}=7$  $\frac{49}{8}=6$  $\frac{43}{8}=5$  $\frac{38}{8}=4$  4    4    4
$P_1$      $P_2$      $P_3$      $P_4$      $P_5$      $P_6$  $P_7$  $P_8$

Execution time = (8+4) × 2.0 = 24 time units.

**CHALMERS**

| Anonymous code | Points for question (to be filled in by teacher) | Consecutive page no. Löpande sid nr  4 |
| Anonym kod   PST431-0007-ADT | Poäng på uppgiften (ifylles av lärare) | Question no. Uppgift nr  2 |

2b) We can ~~tollasce~~ collesce (collapse) the nested loop
   to be of depth 1.

```
#define N 64.
for ( int j=0; j < N×N ; j++) {
    int k = j % N;
    ~~ont~~
    int i = j/N ;
    a[i*N+k] = (i+1)*N+k;
}
```

③

2c)  Overhead > speedup.

$$0.5 > ceil\left(\frac{x}{p}\right) \times 2.0 \text{ where } x \text{ is no. of iterations.}$$

$$\frac{1}{4} > ceil\left(\frac{4096}{p}\right)$$

$$p > 1024.$$

②

After $p=1024$, adding new processors
result in a slowdown.

CHALMERS

Anonymous code
Anonym kod

DIT431-0007-ADT

Points for question
(to be filled in by teacher)
Poäng på uppgiften
(ifylles av lärare)

Consecutive page no.
Löpande sid nr  5

Question no.
Uppgift nr  3

5

3a) $AI = \dfrac{FLOPs}{Bytes} = \dfrac{2}{48 \times 4} = \dfrac{1}{8}$ Flops/Byte, since C has to be loaded and stored, and each various float is 4 bytes.

Single core : performance = $\min\left(4, \dfrac{1}{48} \times 16\right)$

$= 2$ GFLOPs.

execution time $= \dfrac{100 \times 100 \times 100 \times 2}{10^6 \times 2} = 1$ s.

8 cores : performance $= \min\left(8 \times 4, \dfrac{1}{8} \times 16\right)$

$= 2$ GFLOPs.

execution time $= 1$ s

3b) By multiplying the performance of each core by 4, the AI is not changed thus the bottleneck is still Memory. Even in a 32-core machine, the bottleneck is still the same.

3c) multi-core chip execution time $= 1$s.

upper limit $= \dfrac{1}{f}$ where $f$ is serial fraction

$= \dfrac{1}{\frac{0.5 \times 10^{-3}}{1}} = 2000.$

Using my results, since the single core has the same execution time since it is bounded by memory, executing on more cores will not improve result in more speed-up. Optimising the kernel is very suitable here. Optimising the serial portion is less suitable since the ration of $0.5 \times 10^{-3}$s to $1$s is already quite low.

CHALMERS

Anonymous code

Anonym kod

DIT431-6007-ADT

Points for question
(to be filled in by teacher)

Poäng på uppgiften
(ifylles av lärare)

Consecutive page no.
Löpande sid nr  6

Question no.
Uppgift nr  4

5

4a) No. There is a read after write dependency. For example, when $i=0$ we write to $A[0]$, but when $i = N-1$, we read from $A[0]$.

4b).

```
parallel

#omp ppara

# pragma omp parallel for reducon
for (int i=0 ; i<N; i++) {
    if ( i < N/2 ) {
        A[i] += A[N-1-i];
    }

    # pragma omp barrier;
    if ( i >= N/2 ) {
        A[i] += A[N-1-i];
    }
}.
```

**CHALMERS**

5

Anonymous code

Anonymous kod

Points for question
(to be filled in by teacher)

Poäng på uppgiften
(ifylles av lärare)

Consecutive page no.
Löpande sid nr  7

Question no.
Uppgift nr   4

DIT431-0007-ADT

4c) quad += temp has a dependency and there is a risk that
      stores are lost.

First Method:

  We can make quad += temp an atomic instruction using
  #pragma omp atomic. This way every variable can be
  shared, except for temp which should be private.

Second method:
  We can use make quad += temp a critical region so
  that only 1 thread runs the instruction at a given time.
  We use      #pragma omp critical. using shared
    variables, except for temp which should be private. 4

Third method:
  We can use reduce in our for loop   with quad
  and temp as shared variables.
    # pragma omp parallel for reduce(+: quad) private(temp)

4d) I expect using reduce to be the fastest since
every thread can still be fully run asynchronously,
                        atomic
followed by critical then critical, since critical
requires longest waiting. For critical, compiler can
create a lock for the variable quad. For atomic,
a lock can also be used. For reduce, the compiler
can create different quad variables for each thread.

3

**CHALMERS**

| Anonymous code | Points for question | Consecutive page no. |
| | (to be filled in by teacher) | Löpande sid nr  8 |
| Anonym kod | Poäng på uppgiften | Question no. |
| | (ifylles av lärare) | Uppgift nr |
| DIT 431-0007-ADT | | 5 |

5a)

First way: We can use scatter so that the one processor can allocate the chunks of the matrix to the corresponding processors.

~~int id= MPI get MPI id;~~

```
MPI - Scatter ( send buf , bytesToSend, SendType , recvbuf,
        recvCount, recvtype, 0, MPI _ COMMS _ WORLD );
```

\# Now the processors portion of the matrix is in recvbuf

Second way: We can use the root to send each processors portion of the matrix.

```
if    ( get_ MPI_ rank() == 0 )    {
      for  (int i=0; i < P ; i++)    {
      MPI_Send
            MPI_ Send (buf [i] , count, datatype, i, 0, MPI_COMMS_WORLD)
      }
} else    {
                    recvbuf                    get_MPI_rank()
      MPI_ Recv ( buf, count, datatype, dest         , 0, MPI_comm
                                                          world
                                          status );
}
```

**CHALMERS**

| Anonymous code | Points for question (to be filled in by teacher) | Consecutive page no. Löpande sid nr  9 |
| Anonym kod | Poäng på uppgiften (ifylles av lärare) | Question no. Uppgift nr  5 |
| DIT431-0007-ADT | | |

5

5b)

```
int* block = .... // attained from 5a)
int zeros = 0;
for (int i=0; i < N/P; i++) {
    for (int j=0; j<N; j++) {
        if (block [i][j] == 0) {
            zeros++;
        }
    }
}

int* receive;
    :

MPI_Reduce ( &zeros, receive, P, MPI_INTEGER,
MPI_ADDITION, 0, MPI_COMMS_WORLD
);

MPI_Bcast (receive, P, MPI_INTEGER, 0, MPI_COMMS_WORLD
);
```

5c) SSend and SSrecv are synchronous and non blocking, which means SSend only returns when the other side has started receiving and SSrecv only returns when it has ~~started receiving~~ received into the buffer. This means for every SSend, it is able to synchronise with a SSrecv and vice versa, since no deadlocks were observed. Thus, even if they are replaced to non-blocking and wait calls, no deadlock would occur.

2

CHALMERS

Anonymous code

Anonym kod

DIT431-0007-ADT

Points for question
(to be filled in by teacher)

Poäng på uppgiften
(ifylles av lärare)

Consecutive page no.
Löpande sid nr  11

Question no.
Uppgift nr  6

5

6a)

A: & d_row

B: & d_col

C: & d_data

D: & d_x

E: & d_y

F: d_row, &row

G: d_col, & col

H: d_data, & data

I: d_x, & x.

J: 1, N.

K: d_row, d_col, d_data, d_x, d_y, N

L: & y, d_y

6b) ~~Row and colum~~
~~Data and~~
data and col vectors are accessed sequentially in each thread.

-- global -- ~~void~~ spmv (.....)
-- shared -- shared-

CHALMERS

Anonymous code

Anonym kod

DIT431- 0007-ADT

Points for question
(to be filled in by teacher)

Poäng på uppgiften
(ifylles av lärare)

Consecutive page no.
Löpande sid nr  12

Question no.
Uppgift nr  6

5

6b) $x$ offers temporal locality as each value in $x$ is accessed multiple times. Thus we can bring $x$ into shared memory.

```
--global-- void spmv ( .... )      {
      --shared-- x_shared [N];
      int tid = ...
      if ( tid <N)   {
            float dot =0.0;
            for  (int i = row[tid] .... )    {
                  int col_ind = col[i];
                  int~~~~~~~~~
                  if  (x_shared [col_ind] == -1 ) {      // (Some way
                        x_shared [col_ind] = x[col_ind];     of checking
                                                          value has not
                  }                                        been added)
                  dot += data[i] * x_shared [col_ind];
            }
            y[tid]= dot;
      }
}
```

2

**CHALMERS**

Anonymous code

Anonym kod

DIT431 – 0007 – ADT

Points for question
(to be filled in by teacher)

Poäng på uppgiften
(ifylles av lärare)

Consecutive page no.
Löpande sid nr  13

Question no.
Uppgift nr  6

5

6c)

6d) Make use of loop tiling so that we access $x$ in blocks, to and each thread block is able to share the block of $x$ that they require.