Görkem YAR
27970

## Q1:

### Graph (a)

1) Maximum degree of concurrency : 8
2) Critical path length : 4 nodes 3 length
3) Maximum Speedup: $\frac{15}{4} = 3.75$
4) Minimum number of processors for max speed up = 8
5) Maximum Speed up with 4 threads = $\frac{15}{5} = 3$

### Graph (b)

1) Maximum degree of concurrency = 8
2) Critical path length = 4 nodes 3 length
3) Maximum speed up = $\frac{15}{4} = 3.75$
4) Minimum number of processor for max speed up = 8
5) Maximum speed up with 4-threads = $\frac{15}{5} = 3$

### Graph (c)

1) Maximum degree of concurrency: 8
2) Critical Path Length = 7 nodes 6 length
3) Maximum Speed-up = $\frac{14}{7} = 2$
4) Minimum number of processors for max speed up = 3
5) Maximum speed up with 4 threads = $\frac{14}{7} = 2$

### Graph (d)

1) Maximum degree of concurrency: 8
2) Critical path length = 8 nodes 7 length
3) Maximum Speed up = $\frac{15}{8} = 1.75$
4) Minimum Number of processor for max speed up = 2
5) Maximum Speed up with 4 threads = $\frac{15}{8} = 1.75$

## Q2:

<u>revised code</u>

```
int temp;
for (j=1; j<n; j++){
    temp = C[j]
    for(i=1; i<n; i++){
        C[j][i] = C[j-1][i-1] + temp
```

\* improvement-1: change loop order to improve spatial locality. In the previous version the code access the matrix column by column. Now the access is done by row by row. As a result, the cache miss rate will reduce. (Cache Locality) ↗

\* improvement-2: Using a temp variable to store C[j]. Instead of accessing it from memory multiple times, we put it in a register. This reduced memory accesses. (Temporal Locality) ← (Locality in registers) ↗

**Q3:** Using Amdahl's law:

Serial fraction $1-r$; parallel fraction $r$; $T_s$ represents serial time

$T_P(P) = (1-r) \cdot T_s + \dfrac{r \cdot T_s}{P}$   $T_P$ represents parallel time

Amdahl's law for speedup:

$\subseteq T_s / T_P(P)$

$= \dfrac{T_s}{(1-r) \cdot T_s + \dfrac{r \cdot T_s}{P}} = \dfrac{1}{(1-r) + \dfrac{r}{P}}$

if we make $\dfrac{r}{P}$ to small compared to $1-r$ then it can be negligible.

$\lim\limits_{P \to \infty} \dfrac{1}{(1-r) + \dfrac{r}{P}} = \dfrac{1}{1-r}$ // speedup maximum

**Q4:**

a) Latency: The time it takes for the data to travel from memory to it's destination. It is the measurement of delay and unit of latency is seconds.

Bandwidth: The maximum rate for the data to transfer from memory. It can be measured as bit/seconds or byt/seconds. →capacity

*A 'low latency means faster response, high bandwidth means more data can be transmitted at once.

b) Spatial Locality: It means use of data elements within relatively close in memory (Location). If a particular memory accessed nearby locations will be accessed as well.

Temporal Locality: Refers to the reuse of specific data within a small time interval.

Görkem YAR
27970

**Q5:**

$10,6$ Teraflops $= 10,6 \cdot 10^{12}$ floating point operations per second

$720$ 6 Bytes $= \dfrac{720}{4} \cdot 10^9 = 180 \cdot 10^9$ floats per second from global memory

$\searrow$ sizeof (float)

$$\dfrac{10,6 \cdot 10^{12}}{180 \cdot 10^9} = 58,89 \text{ floating point operations per float.}$$

on average 59 operations.

**Q6:**

a) $50.000 / 2048 = 24, 414 \rightarrow \boxed{25}$ blocks    each block has 1024 threads

b) Each warp in CUDA is 32 threads each block has 1024 threads

So each block has $1024/32 = \underline{32}$ warps

total warps $= 25 \cdot 32 = 800$ warps

c) $1024 \times 25 = 25600$ threads in total.

d) Let's consider the following threads with block numbers:

blockIdx.x $= 24$

blockDim.x $= 1024$

id of the last blocks first thread $\rightarrow 49.152$

$50000 - 49152 = 848$ is the threadIdx of the first thread that does not do any job.

block $= 25$ (last block)    warp $= \boxed{848 / 32} = 26$ is the warp number

e) The control diverges for the threads that higher idx or warp number in line 5 and 7.
threads — warps
that are 26 or higher

Also, the first 848 threads of the last block diverge in line 7 as well.
first 28 warps
less than 26

**Q7:**

**a)** Using Amdahl's law: The unparalelisable part of a program determines the maximum speed up.

Serial      Fraction      40 %
parallel    Fraction      60 %

$$T_P(P) = \%40 \cdot T_S + \frac{\%60 \cdot T_S}{P}$$

$$\text{Speed up} = \frac{T_S}{\frac{T_S \cdot 2}{5} + \frac{3}{5} \frac{T_S}{P}} = \frac{1}{\frac{2}{5} + \frac{3}{5 \cdot P}} \qquad \text{if } \lim_{P \to \infty} = \frac{5}{2} = 2.5 \text{ times speed up}$$

to much cores, parallel work is $P$. However, in practice

**b)** theoretically maximum speed up with $P$ processsor can exceed $P$ (we call this super linearity).

the maximum speed up can exceed $P$ (we call this super linearity).

examples
- N queers problem : parallel version can have less work since the solution might appear in one parallel block much quicker than doing all serial work.
- Also, the parallel executions can improve cache-hit ratios this is called resource-based superlinearity.

**c)**
$$E = \frac{S}{P} = \frac{T_S}{P \cdot T_P} = \frac{n}{P \cdot (n/p + \log p)} = \frac{n}{n + p \log p}$$

**d)** $E = \Theta(1)$ to be cost optimal. The cost optimality can also be written in the form

$$\overbrace{P \cdot (n/p + \log p)}^{} = \underbrace{P \cdot T_P = \Theta(T_S)}_{}$$

$$P(n/p + \log p) = \Theta(n)$$

$$= n + p \log p = \Theta(n)$$

$$\Rightarrow p \log p = \Theta(n)$$

$$\Rightarrow p \log p \leq kn$$

$$\frac{n}{\log n}(\log n - \log \log n) \to \qquad \boxed{P = \frac{n}{\log n}}\ \text{by substitution}$$

$$n - \frac{\log \log n}{\log n} \leq k \cdot n \quad \text{holds}$$

$$\text{to smit}\ \frac{}{\log n}$$