

Aabash

2018B4A70887P

Design:

a) Huffman algo is not parallelized because

→ Based on literatures available online, the overhead of parallelizing is high.

→ main reason is that our input space is from sparse space - printable ascii characters. Also, ~~since~~ if this assumption is not true, even then the ~~new~~ symbols are bounded by 256 because splitting is done based on the bytes.For files  $\geq 10$  mb, the time for running Huffman, given the frequency of each character, is very low.

b) For encoding,

Time complexity  $= O\left(\frac{n}{P}\right)$ ,

i) File is opened using MPI-open and each process gets its own copy of data, which is distinct and complete data is separated b/w processes. Data is divided on byte level and ~~for start index~~ every process gets no. of bytes  $\left\{ \left\lfloor \frac{\text{size}}{P} \right\rfloor, \left\lceil \frac{\text{size}}{P} \right\rceil \right\}$

→ Almost uniformly distributed for large files.

ii) Each process calculates frequency of each of byte sequences. →  $O\left(\frac{n}{P}\right)$  time

iii) we need cumulative frequency, so we do All reduce on vectors of size 256 across all processes. Now, size of this vector  $\leq 256$   
 ⇒ overhead is  $O(\log P)$

iv) All processes will get total frequency and use them to construct Huffman tree →  $\leq 256$  characters  
 →  $O(1)$  time.

(This is repeated calculation → serial part this is not parallelized)

(v) Encoding is done by simple linear traversal  
 $\rightarrow O(\frac{n}{p})$  time

(vi) File is stored  $\rightarrow O(\frac{n}{p})$  time

vii) File offsets  $\rightarrow$  MPI Scan  $\rightarrow O(\log p)$

$\Rightarrow$  Overall time complexity =  $O(\frac{n}{p})$

$\hookrightarrow$  Serial time is only for construction of Huffman tree (because calculation is repeated).  
while rest part is ~~for~~ parallelized completely.

Since overhead of transmission is small, it is ignored in calculations. ( $n \gg p$ ) and msg size is fixed

(c) Similar analysis can be done for decoding to get  $O(\frac{n}{p} + 2 \log p)$  as total time.

Part b: To calculate the value of  $f$ , following method is used:

①  $T_s'$  &  $T_p'$  are defined.

$\rightarrow$  Since only serial part, effectively, is construction of Huffman tree,  $T_s'$  is user + sys time for construction of tree and rest is in  $T_p'$

$\rightarrow$  Calculated in code

② let  $\sigma(n)$  be <sup>time for</sup> sequential part of code  
 $\phi(n)$   $\rightarrow$  parallel

when executed sequentially

$$\Rightarrow T_s' = \sigma(n)$$

$$T_p' = \frac{\phi(n)}{p}$$

③ By definition,  $T_1 = \text{time to run code sequentially}$   
 $= \sigma(n) + \phi(n)$   
 $= T_1' + p T_p'$

$T_p = \text{time to run code parallelly}$   
 $= \sigma(n) + \frac{\phi(n)}{p}$   
 $= T_1' + T_p'$

④ By definition,

$$\text{efficiency} = \frac{T_1}{T_p} \leftarrow \text{calculate in code.}$$

⑤ for fraction of code running sequentially,

$$\text{efficiency} = \frac{T_1}{T_p} = \frac{1}{f + \frac{(1-f)}{p}} = \frac{T_1' + p T_p'}{T_1' + T_p'}$$

$$\Rightarrow \frac{T_1' + p T_p'}{T_1' + T_p'} = \frac{1}{f + \frac{(1-f)}{p}}$$

$$\Rightarrow \frac{\sigma(n) + \phi(n)}{\sigma(n) + \frac{\phi(n)}{p}} = \frac{1}{f + \frac{(1-f)}{p}}$$

Comparing,  $f = \frac{\sigma(n)}{\sigma(n) + \phi(n)} = \frac{T_1'}{T_1' + p T_p'} = \frac{T_1'}{T_1}$

for encoding a file of 1 MB,

$$\text{efficiency (avg)} = 0.003$$

for decoding, avg efficiency = 0.004



Part c

Time complexity  $\rightarrow$  derived already as  $O\left(\frac{n}{p}\right)$

for both.

$$\text{speed up} = \frac{T_1}{T_p} = \frac{O(n)}{O\left(\frac{n}{p}\right)} \sim O(p)$$

$$\text{Efficiency} = \frac{\text{Speedup}}{p} = \frac{O(p)}{p} \sim 1$$

$\rightarrow$  almost no overhead for large input

$$\text{Cost} = p T_p = p \cdot O\left(\frac{n}{p}\right) = O(n)$$

$\rightarrow$  This is optimal

Ans (c)

Time complexity  $\rightarrow$  derived already  $= O\left(\frac{n}{p} + 2 \log p\right)$

$$\text{speed up} = \frac{T_1}{T_p} = \frac{n}{\frac{n}{p} + 2 \log p}$$

$$\text{Efficiency} = \frac{\text{Speedup}}{p} = \frac{n}{n + 2p \log p}$$

$$\text{Cost} = p T_p = p \left( \frac{n}{p} + 2 \log p \right) = n + 2p \log p$$

for  $n \gg p$ ,  $\sim O(n)$

for isoefficiency,

$$E_1 = \frac{n_1}{n_1 + 2p_1 \log p_1}$$

$$E_2 = \frac{n_2}{n_2 + 2p_2 \log p_2}$$

For iso efficiency,

$$E = \frac{1}{1 + \frac{2Kp \log p}{n}}$$

$$\Rightarrow 1 + \frac{2Kp \log p}{n} = \frac{1}{E}$$

$$\Rightarrow \frac{2Kp \log p}{n} = \frac{1}{E} - 1$$

$$\Rightarrow n = \frac{2Kp \log p}{\frac{1}{E} - 1} \Rightarrow n \propto 2Kp \log p$$

$\Rightarrow$  ISO Efficiency function:  $w = 2Kp \log p$

(d) optimal no. of processors:

$$T_p = \Theta\left(\frac{w}{p}\right) = \Omega\left(\frac{w}{f^{-1}(w)}\right)$$

where  $w = \Theta(f(p))$

Here,  $w = f(p) = 2Kp \log p$

$$\begin{aligned} \Rightarrow p &\approx \frac{w}{\log w} \Rightarrow T_p = \Omega\left(\frac{w \log w}{w}\right) \\ &= \Omega(\log w) \end{aligned}$$