

Parallel implementation of Huffman Code using native C++ threads and FastFlow library

Parallel and Distributed Systems: Paradigms and Models

University of Pisa

Project Report

Matteo Tolloso

`m.tolloso@studenti.unipi.it`

Roll number: 598067

September 1, 2023

1 Introduction

The Huffman code is an efficient lossless compression code based on the probability of each character.

To build the optimal code for a specific text we have to:

1. count the number of occurrences of each character in the text;
2. build the binary tree that represents the code;
3. encode the file.

2 Overview

We are facing a problem that can be divided in three stages. In particular it is a *data parallel* task, since we have all input available at the beginning of the computation.

2.1 Counting the number of occurrences

As stated above, the first stage is a counting one. The asymptotic sequential complexity of this part is $\theta(m)$ where m is the number of characters in the file. From a parallel/distributer point of view this is clearly a *map-reduce* operation.

Map The *Map* part can be execute in parallel dividing the file into chunks, the workers count the occurrences in a chunk of the file. This operation has to deal with the disk. If we consider the reading of the disk as a sequential operation things became more difficult because it's no longer a data parallel problem but a stream parallel one. In this setting we can describe the process as `pipe(reading, farm(counting, nw))`, the completion time of this process is the time needed to read the file from the disk, under the assumption that the farm has the right number of workers to match the speed of the disk. This approach is the one that minimize both the completion time and the number of workers but it causes a lot of comunication overhead, needs some tuning of the chunksize to send and of the scheduler's policy and is in general more complex to implement.

If we instead consider the reading of the disk as a datata parallel operation that consists in moving data from the disk to the main memory, we can use the *Map Fusion* theorem and transform the program in `map(read-count, nw)`. This solution minimize the communication overhead, the completion time and the complexity of the implementation.

Furthermore, the tests that i did mapping the file in main memory and reading it with multiple threads, showed that the parallelization also improves the performance of the read operation. This is probably due to how the SSD works and the caching systems.

Reduce After the *Map* operation we end up with a number of counts vectors equal to the number of chunks the file was divided into (that in our case is equal to the number of workers). The *Reduce* operation is again a parallel one, this time each reducer takes a subset of the alphabet and sums the occurrences of each character in that subset. It's useless to have a number of workers greater than the number of different characters in the file.

2.2 Building the binary tree

The second stage is the building of the binary tree. This is a more difficult operation to parallelize since most of the operations are sequential. Furthermore, the complexity of this stage is $\theta(A \times \log(A))$ where A is the number of different symbols (128), so basically it is a constant in our case. Tests showed that the time needed for this stage is completely negligible with respect to the other stages.

2.3 Encoding the file

The last stage is the *encoding* of the file. This is a *Map* operation, since each character has to be replaced with its code and written on the disk. We can make a similar reasoning as the one made for the counting stage about the *Map Fusion* theorem.

Unfortunately, the length of the final text can only be known after each character has been encoded (because the encoding of each character has a different length), so the actual writing needs a step of synchronization. I solved this problem dividing this stage in thread parts:

1. *encoding*: each worker encodes a chunk of the file.
2. *balancing*: the encoded chunks sizes are made multiple of 8 and the index where the writing should start is computed. This is a sequential synchronization step but the time needed is negligible.
3. *compressing and writing*: each worker takes a chunk of the encoded file and writes it on the disk grouping the bits in bytes.

It's fundamental to notice that the *balancing* step makes the encoded chunks independent one from another, so the *compressing and writing* can be a parallel operation.

3 Technical details

3.1 Implementation

The FastFlow implementation relies on the application API `ParallelForReduce` that manages the work for each thread and the synchronizations.

In the thread implementation *nw* threads are spawned and a shared queue is used to send them the functions to compute. Each thread calculates its part of the work and after the computation they wait on an explicit barrier. When all the threads have finished the barrier is opened and the next step is executed. The termination is managed with the `optional` type.

The sequential implementation is easy, the characters are iterated one by one while counting the occurrences and after the creation of the code the file is encoded and written.

In all three cases the file is mapped in memory and read as an array of characters. Likewise, the encoded file is mapped in the main memory and written. At the end you make sure that it is synchronized with the disk.

In the parallel implementations there are *nw* count vectors, each one is allocated by each different worker. In the reduce part each worker creates a `map` data structure that are then joined together with a sequential operation. FastFlow automatically manages the allocations in this phase with the `ParallelForReduce` class. In the *encode* phase each worker allocates a tuple containing a `deque` data structure and a `int`. The `deque` is used to store the encoded bits and the `int` will be used in the balancing stage to store the byte to start writing in the encoded file. The balancing does precisely this: moves bits between the `deques` in order to make them multiple of 8 (padding the last one) and computes the starting byte. This operation takes only few usec thanks to the `deque` data structure. After this stage of *balancing* we end up with *nw* `deque` multiple of 8 and for each one we know where it has to be written in the file. The last stage, *compressing and writing*, can now start in parallel. Each worker takes a `deque` and writes it on the file grouping the bits in bytes.

3.2 Overheads

False Sharing The false sharing problem is avoided since each worker writes on a completely different array: the counting arrays and the chunk-encoding arrays are allocated by each worker.

Heap pressure The access to the heap is mutual exclusive, so a high number of allocation reallocation can cause a big overhead. The use of the jemalloc library helps a lot in this case. Details are discussed in section 4.

Load balancing The parallel implementations use a static scheduling. During the *counting* operation the file is equally divided between the workers i.e. each worker counts the same number of characters. In the *reduce* phase each worker takes an equal subset of characters and sums the occurrences. In both phases could happen that a worker has to deal with a bigger number with respect to others (because of different character frequency), but there are only *+1* operations and the time should not depend on the size of the number. In the *encoding* phase each worker takes a chunk to encode. This

part can be really unbalanced if the original file has somewhere a lot of alligned equal character, infact, this character will probability have a short code and the worker that encodes that chunk has to do fewer memory reallocations.

These observations could prompt us to use dynamic scheduling. The FastFlow implementation easily allow us to do that changing only one parameter, however some modifications have to be done in the data structures that collect the results and in general the sequential fraction of the program increases due to the operations to putting data back in order. For a fair comparison I decided to use static scheduling in both implementations.

Synchronization In the FastFlow implementation the synchronization is completely managed by the library. One set of threads is spown at the beginning and the runtime support manages the queues and the implicit barriers. In the native threads implementation I had to manually managed the synchronization. The easiest way would have been to spown and join a set of threads for each stage, each time with the assigned funcion and arguments. This approach would have been really simple but it would have caused a lot of overheads since from some tests on the reference machine the creation and join of a thread takes about $70\mu s$, while the insertion of a task in a shaerd queue takes about $1\mu s$ (and the creation of the shared queue takes $4\mu s$). So the the latter is the method I used.

4 Tests

The table 1 shows the time of the various stages of the sequential implementation. The great part of the time is spent on the *encoding* and *compressing and writing* phases.

Stage	Time
read and count	26 s (25928430 μs)
huffman	0.000102 s (103 μs)
encoding	212 s (212385536 μs)
compressing and writing	318 s (317802196 μs)
Total	567 s (566897566 μs)

Table 1: Sequential times for 8GB file of random characters. Averaged over 10 runs.

In tables 2 and 3 we can see some measures of the FastFlow implementation and the native threads one. The total do not exactly corrspond to the sum of the single stages because they didn't take into account the initialization of the memory and other structures. The measureas of the single stages refer only to the actual computation while the "total" refers to the time from the start of the program to the end.

Stage	Time	Speedup	Efficiency
read and count	0.8 s (839930 μs)	30.87 x	0.96
huffman	0.000077 s (77 μs)		
encoding	9.7 s (9725888 μs)	21.83 x	0.68
balancing	0.000073 s (73 μs)		
compressing and writing	20.8 s (20835071 μs)	15.25 x	0.48
Total	36.3 s (36295647 μs)	15.61 x	0.49

Table 2: Parallel times with FastFlow implementation for 8GB file of random characters. 32 physical core machine. Averaged over 10 runs.

Stage	Time	Speedup	Efficiency
read and count	0.9 s (894301 μs)	28.99 x	0.90
huffman	0.000095 s (95 μs)		
encoding	10.3 s (10269624 μs)	20.68 x	0.65
balancing	0.000027 s (27 μs)		
compressing and writing	21.2 s (21213486 μs)	13.86 x	0.43
Total	33.3 s (33328868 μs)	17.00 x	0.53

Table 3: Parallel times with native threads implementation for 8GB file of random characters. 32 physical core machine. Averaged over 10 runs.

The *read and count* stage has a great speedup, in particular the application API of FastFlow allows an almost linear speedup in this operation and more in general the speedup of the actual computation is always better with FastFlow than with the native threads implementation. If we instead consider the total speedup, threads are slightly better probably due to less overhead in the management.

The *compress and write* stage has the worst speedup because it involves writing on disk, so it incorporate a substantial sequential fraction.

The *encoding* phase is independet from the disk and I expected a better speedup. Probably the overhead is caused by the memory reallocations needed to store the encoded characters, ence a competition to access the heap despite the use of the jemalloc library. Even allocate a lot of memory at the beginnig is not a solution because the threads will compete the same. One possible solution could be switch to arena mode immediately but it is ouside my control.

In the figure 1 the completion time of the two parallel implementation is compared. The FastFlow result is quite surprising because it achieves the minimum time with only 16 threads, and starts to worsen the performance increasing the parallel degree. A similar behaviour can be observed with the native thread, wehere is not worth to increase the

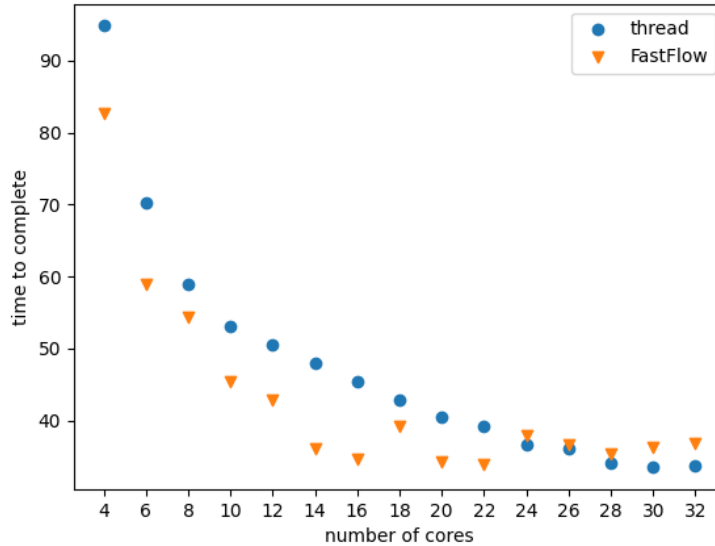


Figure 1: Time for encoding of a 8GB file of random characters, in seconds. Averaged over 10 runs.

number of threads over 28. The same results can be seen under a different prospective in the efficiency plot, figure 2. I suspect that the great efficiency of FastFlow with 16 cores is correlated with the architecture of the machine used for the experiments (2×16 cores) and some optimizations made by FastFlow and not present with the native threads implementation (for example thread pinning).

From the figure 2 we can also observe that the efficiency is greater than one. This can happen when the speedup is greater than the parallel degree because the original problem became much easier when it is divided into subproblems. It might be that allocate and work a single 8 GB block is much slower than allocate 32 different 256 MB blocks, especially with the optimizations made by the jemalloc library and FastFlow. For example could happen that the sequential process allocates the memory on the RAM bank attached to the current core and then is moved to another core or even another socket causing a lot of delay in accessing memory; in the parallel process is more likely that each thread allocate the memory on the nearest RAM bank and then remains pinned to the same core, thanks to systems like the “thread pinning”, that is present in FastFlow.

These are only hypothesis difficult to verify, however the test reported in table 4 shows how important is the use of the jemalloc library for these implementations, especially with FastFlow, that suffers by an order of magnitude the unoptimized standard library. This could partially explain the superlinear speedup observed before.

In table 5 the same test of table 2 is repeated with 16 cores to confirm that the best efficiency with FastFlow is obtained with 16 threads. Results show that the efficiency is

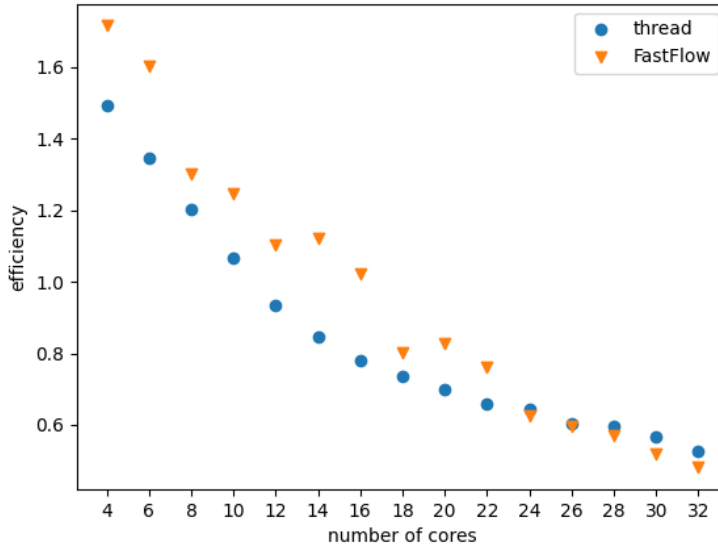


Figure 2: Efficiency in the encoding of a 8GB file of random characters. Averaged over 10 runs.

surprisingly high since the speedup is almost linear. The 0.9 efficiency in *compressing and writing* stage shows that we are near the bandwidth of the disk and the 1.06 efficiency in the *encoding* stage is probably due to hypothesis above.

	malloc	jemalloc	speedup
threads	196 s (169339515 μs)	33.3 s (33328868 μs)	5.08 x
FastFlow	415 s (414727131 μs)	36.2 s (36295647 μs)	11.42 x

Table 4: Time to encode a 8GB file of random characters. Standart malloc library vs jemalloc. 32 physical core machine. Averaged over 10 runs. The term “speedup” here has a different meaning than the usual one: it is the ratio between the time with the standart library and the time with jemalloc.

It is useful to know when it is worth to use a parallel implementation instead of a sequential one. Give a precise answer to this question is difficult because the variables involved are many (file size, single core speed, parallel degree, etc.). However, in figure 3 we can see the time to encode different sized files fixing the number of workers to 32. We can see that with half a milion characters the parallel implementations are already faster than the sequential one. After the consideration made above about the efficiency, we can imagine that reducing a bit the number of workers reduces the overhead without

Stage	Time	Speedup	Efficiency
read and count	1.6 s (1633635 μ s)	15.87 x	0.99
huffman	0.000095 s (95 μ s)		
encoding	12.5 s (12521952 μ s)	16.96 x	1.06
balancing	0.000027 s (27 μ s)		
compressing and writing	22 s (22056740 μ s)	14.04 x	0.90
Total	38.7 s (38714097 μ s)	14.64 x	0.92

Table 5: Parallel times with FastFlow implementation for 8GB file of random characters. 16 physical core machine. Averaged over 10 runs.

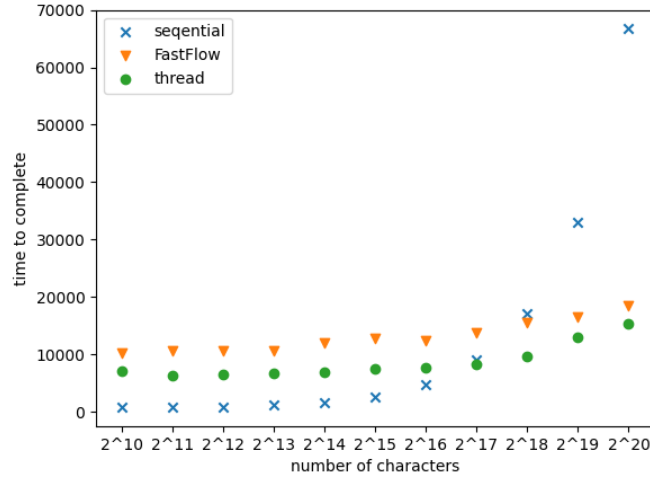


Figure 3: Time for encoding different sized files of random characters, in μ s, with 32 physical cores. Averaged over 10 runs.

increasing the completion time, so it would become useful to use parallel implementations even with smaller files.

5 Conclusions

The tests showed that predict the behaviour of a parallel application is really difficult, as well as write a flexible software capable to adapt to the underlying architecture that could strongly influence the performance.

FastFlow has a lot of options and optimizations and the code results more compact. However, for the same reason it is more difficult to predict the behaviour of the application and to tune the parameters; in general the tests showed a really good performance.

The native threads implementation code is longer and built around this use case. Implement more complex policy and optimizations would be very difficult and error-prone. The advantage is the almost absent runtime overhead that brings to a slightly better completion time.