PARALLEL DATA COMPRESSION WITH BZIP2

Jeff Gilchrist Elytra Enterpises Inc. Unit 1, 5370 Canotek Road Ottawa, Canada K1J 9E6 jeff@gilchrist.ca

ABSTRACT

A parallel implementation of the bzip2 block-sorting lossless compression program is described. The performance of the parallel implementation is compared to the sequential bzip2 program running on various shared memory parallel architectures.

The parallel bzip2 algorithm works by taking the blocks of input data and running them through the Burrows-Wheeler Transform (BWT) simultaneously on multiple processors using pthreads. The output of the algorithm is fully compatible with the sequential version of bzip2 which is in wide use today.

The results show that a significant, near-linear speedup is achieved by using the parallel bzip2 program on systems with multiple processors. This will greatly reduce the time it takes to compress large amounts of data while remaining fully compatible with the sequential version of bzip2.

KEY WORDS

Parallel Computing, Parallel Algorithms, Data Compression, Software Tools

1 Introduction

Parallel computing allows software to run faster by using multiple processors to perform several computations simultaneously whenever possible. The majority of general purpose software is designed for sequential machines and does not take advantage of parallel computing. Lossless data compression software like bzip2, gzip, and zip are designed for sequential machines. The data compression algorithms require a great deal of processing power to analyze and then encode data into a smaller form. Creating a general purpose compressor that can take advantage of parallel computing should greatly reduce the amount of time it requires to compress files, especially large ones.

This paper explores several ways of implementing a parallel version of the Burrows-Wheeler Transform (BWT) and describes one implementation that modifies the popular bzip2 compression program. The limited research papers available on parallel lossless data compression use theoretical algorithms designed to run on impractical parallel machines that do not exist or are very rare in the real world [10, 11]. The purpose of a parallel bzip2 compression pro-

gram is to make a practical parallel compression utility that works with real machines. The bzip2 program was chosen because it achieves very good compression and is available in library and source form, free under an open source license.

One method of parallelizing the BWT algorithm has been implemented and tested on several shared memory parallel machines such as a dual processor AMD Athlon MP and Itanium2 system, and a 24 processor SunFire 6800 system. The results show that near-linear speedup is achieved.

The remainder of the paper is organized as follows. In Section 2, the relevant literature on data compression is reviewed. Section 3, shows several possible ways to parallelize the BWT algorithm. Section 4, presents the parallel bzip2 algorithm. Section 5 gives the performance results using the algorithm on several parallel architectures. Section 6 concludes the paper.

2 Background

Lossless data compression is used to compact files or data into a smaller form. It is often used to package up software before it is sent over the Internet or downloaded from a web site to reduce the amount of time and bandwidth required to transmit the data. Lossless data compression has the constraint that when data is uncompressed, it must be identical to the original data that was compressed. Graphics, audio, and video compression such as JPEG, MP3, and MPEG on the other hand use lossy compression schemes which throw away some of the original data to compress the files even further. We will be focusing only on the lossless kind. There are generally two classes of lossless compressors: dictionary compressors and statistical compressors. Dictionary compressors (such as Lempel-Ziv based algorithms) build dictionaries of strings and replace entire groups of symbols. The statistical compressors develop models of the statistics of the input data and use those models to control the final output [5].

2.1 Dictionary Algorithms

In 1977, Jacob Ziv and Abraham Lempel created the first popular universal compression algorithm for data when no a priori knowledge of the source was available. The LZ77 algorithm (and variants) are still used in many popular compression programs today such as ZIP and GZIP. The compression algorithm is a dictionary based compressor that consists of a rule for parsing strings of symbols from a finite alphabet into substrings whose lengths do not exceed a predetermined size, and a coding scheme which maps these substrings into decipherable code-words of fixed length over the same alphabet [12]. Many variants and improvements to the LZ algorithm were proposed and implemented since its initial release. One such improvement is LZW (Lempel-Ziv-Welch) which is an adaptive technique. As the algorithm runs, a dictionary of the strings which have appeared is updated and maintained. The algorithm is adaptive because it will add new strings to the dictionary [7].

Some research on creating parallel dictionary compression has been performed. Once the dictionary has been created in an LZ based algorithm, a greedy parsing is executed to determine which substrings will be replaced by pointers into the dictionary. The longest match step of the greedy parsing algorithm can be executed in parallel on a number of processors. For a dictionary of size N, 2N-1processors configured as a binary tree can be used to find the longest match in $O(\log_N)$ time. Each leaf processor performs comparisons for a different dictionary entry. The remaining N-1 processors coordinate the results along the tree in $O(\log_N)$ time [10]. Stauffer a year later expands on previous parallel dictionary compression schemes on the PRAM by using a parallel longest fragment first (LFF) algorithm to parse the input instead of the greedy algorithm [11]. The LFF algorithm in general performs better than the greedy algorithm but is not often used in sequential implementations because it requires two passes over the input. With the PRAM, using the LFF parsing can be performed over the entire input in one step. This algorithm is not practical since very few people if anyone would have a machine with enough processors to satisfy the requirements. Splitting data into smaller blocks for parallel compression with LZ77 yields speedup but poor compression performance because of the smaller dictionary sizes. A cooperative method for building a dictionary achieves speedup while keeping similar compression performance [6].

2.2 Statistical Algorithms

Statistical compressors traditionally combine a modeling stage, followed by a coding stage. The model is constructed from the known input and used to facilitate efficient compression in the coder. Good compressors will use a multiple of three basic modeling techniques. Symbol frequency associates expected frequencies with the possible symbols allowing the coder to use shorter codes for the more frequent symbols. Symbol context has the dependency between adjacent symbols of data usually expressed as a Markov model. This gives the probability of a specific symbol occurring being expressed as a function of the previous n symbols. Symbol ranking occurs when a symbol

"predictor" chooses a probable next symbol, which may be accepted or rejected [5].

Arithmetic coding is a statistical compression technique that uses estimates of the probabilities of events to assign code words. Ideally, short code words are assigned to more probable events and longer code words are assigned to less probable events. The arithmetic coder must work together with a modeler that estimates the probabilities of the events in the coding. Arithmetic coding runs more slowly and is more complex to implement than LZ based algorithms [8].

The PPM (Prediction by Partial Match) algorithm is currently the best lossless data compression algorithm for textual data. It was first published in 1984 by Cleary and Witten [3]. PPM is a finite-context statistical modeling technique which combines several fixed-order context models to predict the next character in the input sequence. The prediction probabilities for each context are adaptively updated from frequency counts. The basic premise of PPM is to use the previous bytes in the input stream to predict the following one. A number of variations and improvements have been made since then.

2.3 Burrows-Wheeler Transform

The Burrows-Wheeler transform [1] is a block-sorting, lossless data compression algorithm that works by applying a reversible transformation to a block of input data. The transform does not perform any compression but modifies the data in a way to make it easy to compress with a secondary algorithm such as "move-to-front" coding and then Huffman, or arithmetic coding. The BWT algorithm achieves compression performance within a few percent of statistical compressors but at speeds comparable to the LZ based algorithms. The BWT algorithm does not process data sequentially but takes blocks of data as a single unit. The transformed block contains the same characters as the original block but in a form that is easy to compress by simple algorithms. Same characters in the original block are often grouped together in the transformed block.

The algorithm works by transforming a string S of N characters by forming the N rotations (cyclic shifts) of S, sorting them lexicographically, and extracting the last character of each of the rotations. A string L is then formed from these extracted characters, where the ith character of L is the last character of the ith sorted rotation. The algorithm also computes the index L of the original string L in the sorted list of rotations. With only L and L there is an efficient algorithm to compute the original string L when undoing the transformation for decompression. They also explain that to achieve good compression, a block size of sufficient value must be chosen, at least L0 kilobytes. Increasing the block size also increases the effectiveness of the algorithm at least up to sizes of several megabytes.

The BWT can be seen as a sequence of three stages: the initial sorting stage which permutes the input text so similar contexts are grouped together, the Move-To-Front stage which converts the local symbol groups into a single global structure, and the final compression stage which takes advantage of the transformed data to produce efficient compressed output [5]. Present implementations of BWT require about 9 bytes of memory for each byte of data, plus a constant 700 kilobytes. For example, to compress 1 megabyte of data would require about 10 megabytes of memory using this algorithm.

Before the data in BWT is compressed, it is run through the Move-To-Front coder. The choice of MTF coder is important and can affect the compression rate of the BWT algorithm. The order of sorting determines which contexts are close to each other in the output and thus the sort order and ordering of source alphabet can be important. Many people consider the MTF coder to be fixed, but any reversible transformation can be used for that phase [2].

It was discovered that with sequences of length n taken from a finite-memory source that the performance of BWT based algorithms converge to the optimal performance at a rate of $O(\frac{\log_n}{n})$ surpassing that of LZ77 which is $O(\log_{\frac{\log_n}{\log_n}})$ and LZ78 which is $O(\frac{1}{\log_n})$ [4].

2.4 BZIP2

The BZIP2 program compresses files using the Burrows-Wheeler block-sorting text compression algorithm, and Huffman coding. Compression is generally considerably better than that achieved by more conventional LZ77/LZ78-based compressors, and approaches the performance of the PPM family of statistical compressors while being much faster [9]. BZIP2 is freely available open source software popular in the Unix world. It is used for compressing software distributions and source code such as the Linux kernel among other things.

3 Parallelizing BWT

Three ways to parallelize the BWT algorithm in a shared memory architecture are considered.

3.1 Parallel Sort

Burrows and Wheeler explain that much of the time, the BWT algorithm is sorting [1]. The initial stage of the BWT algorithm creates a matrix from the input data which needs to be sorted lexographically. Implementing a parallel sorting algorithm may increase the speed of the BWT on a parallel machine.

3.2 Multiple Blocks

The BWT algorithm in most implementations requires about 10 times the amount of memory as input data. In order to cope with large amounts of data, the BWT algorithm will naturally split the data into independent blocks of a predetermined size before the data is compressed. The

blocks are then processed through the BWT algorithm. Since each block is independent, it should be possible to run the BWT algorithm on multiple blocks of data simultaneously and achieve speedup on a parallel machine. The separate blocks are then concatenated back together again to form the final compressed file. The sequential version of BWT will process the blocks in order, so does not need to worry about order when writing the output to disk. A parallel implementation will need to keep track of block ordering and write the compressed blocks back to disk in the correct order.

3.3 Combination

Depending on the number of processors on the parallel machine, it may be useful to combine processing multiple blocks of data simultaneously with using a parallel implementation of the lexicographical sorting algorithm. Tests would have to be performed in order to determine if using a combination is better than either of the solutions alone.

4 Parallel BZIP2

The method of processing multiple blocks of data simultaneously with the BWT algorithm was implemented. The popular bzip2 program by Julian Seward [9] uses the BWT algorithm for compression. It was chosen because it is available in library and source form, free under an open source license. This gives the benefit of not having to implement the BWT from scratch and allows the compressed files generated from the parallel version to be compatible. Anyone with the sequential version of bzip2 would be able to uncompress files created with the parallel version and vice-versa. The parallel version of bzip2 has been named pbzip2.

4.1 Sequential bzip2

The bzip2 program processes data in blocks ranging from 100,000 bytes to 900,000 bytes in size, depending on command line switches. The default block size is 900,000 bytes. It reads in data 5,000 bytes at a time, until enough data to fill a block is obtained. It will then process the block with the BWT algorithm and write the compressed output to disk. This continues until the entire set of data is processed; see Figure 1. Since bzip2 is a sequential program, it will only start processing the next block after the previous one has been completed.

4.2 pbzip2

For pbzip2 to achieve speedup on parallel machines, a multi-threaded design was created using pthreads in C++. The code is generic enough that it should compile on any C++ compiler that supports pthreads. The code is linked

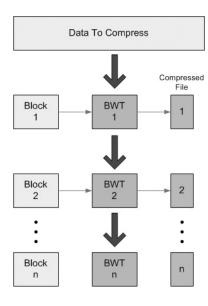


Figure 1. BZIP2 flow diagram

with the libbzip2 library to access the BWT algorithm (http://sources.redhat.com/bzip2/).

The pbzip2 program also works by splitting the data into blocks of equal size, configurable by the user. Instead of reading in 5,000 bytes at a time, pbzip2 reads in an entire block at once which increases performance slightly, even with a single processor. The user can also specify the number of processors for pbzip2 to use during compression. It supports compressing single files just as the bzip2 program does. If multiple files are needed to be compressed together, they should first be processed with TAR to create a single archive, and then the .tar file compressed using pbzip2.

A FIFO (first in, first out) queue system is used in pbzip2 with a producer/consumer model. The size of the queue is set to the number of processors pbzip2 was configured to use. This gives a good balance between speed and amount of memory required to run. Setting the queue size larger than the number of processors did not result in any speedup but significantly increased the amount of memory required to run. If the default block size of 900,000 bytes is used and the number of processors requested was two, pbzip2 will read the first two blocks of 900,000 bytes from the input file and insert the pointers to those buffers into the FIFO queue. Since pbzip2 is using two processors, two consumer threads are created. Mutexes are used to protect the queue so only one thread can read from or modify the queue at a time. As soon as the queue is populated, the consumer threads will remove items from the queue and process the data using the BWT algorithm. Once the data is processed, the memory for the original block is released and the pointer to the compressed block along with the block number is stored in a global table. A file writing thread is also created which will go through the global table and write the compressed blocks to the output file in

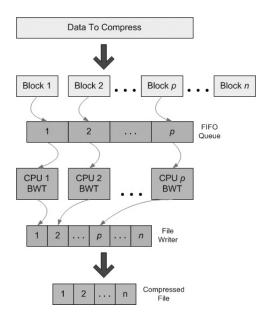


Figure 2. Parallel BZIP2 flow diagram

the correct order. Once free space in the queue is available, more blocks of data will be read from the file and added to the queue. This process continues until all blocks have been read, at which point a global variable is set to notify all the threads that the file has been read in its entirety. The consumer threads will continue to process data blocks until the queue is empty and there are no more blocks going to be added to the queue. The pbzip2 program finishes when the file writing thread has finished writing all the compressed blocks to the output file in their correct order; see Figure 2.

The data being compressed is split into independent parts so minimal communication between processors is required. Also, the BWT algorithm is CPU intensive so distributing the data between processors and the disk I/O for reading and writing should have minimal impact on time. These factors should allow the pbzip2 design to achieve significant, possibly near-linear speedup.

5 Performance Results

The pbzip2 program was compiled on several different parallel platforms using the gcc compiler and the libbzip2 1.0.2 library. The input file used for testing the performance of bzip2 and pbzip2 was a database of elliptic curve distinguished points which is 1,966,717,056 bytes (1.83 GB) uncompressed.

All experiments were measured as wall clock times in seconds. This includes the time taken to read from input files and write to output files. Each experiment was carried out three times and the average of the three results was recorded. Since bzip2 does not take advantage of multiple processors, the running time is the same no matter how many processors the system has.

Due to space constraints, only one set of results are

Athlon-MP (2 x 2.1GHz)

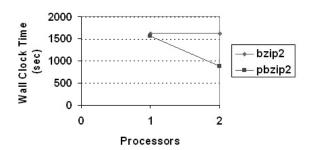


Figure 3. Wall clock time in seconds as a function of the number of processors to compress 1.83 GB on an Athlon-MP 2600+ system.

presented. Multiple experiments were performed and show similar results which are available at http://compression.ca.

5.1 Athlon-MP 2600+

The platform used in this experiment was an AMD Athlon-MP 2600+ machine with two 2.1 GHz processors, 256KB L2 cache, and 1 GB of RAM. The software was compiled with gcc v3.3.1 under cygwin 1.5.10-3 on Windows XP Pro.

The pbzip2 program was run once with the -p1 switch for 1 processor, and tested again with the -p2 switch for 2 processors.

The bzip2 program establishes a baseline for performance. Running pbzip2 on one processor showed that it was slightly faster than bzip2. This can probably be attributed to the fact that pbzip2 will read larger chunks of data at a time than bzip2, reducing the number of reads required. It was observed that when running pbzip2 on two processors, 87.9% of linear speedup was achieved. A graph of the wall clock time as a function of the number of processors is shown in Figure 3.

5.2 Itanium2

The platform used in this experiment was an Intel Itanium2 machine with two 900 MHz processors, 1.5MB L3 cache, and 4 GB of RAM. The software was compiled with gcc v2.96 and run on RedHat Linux 7.2 with the 2.4.19-4 SMP kernel.

The pbzip2 program was run once with the -p1 switch for 1 processor, and tested again with the -p2 switch for 2 processors.

The bzip2 program establishes a baseline for performance. Running pbzip2 on one processor showed that it was slightly faster than bzip2. This can probably be attributed to the fact that pbzip2 will read larger chunks of data at a time than bzip2, reducing the number of reads required. It was observed that when running pbzip2 on two

Itanium 2 (2 x 900 MHz)

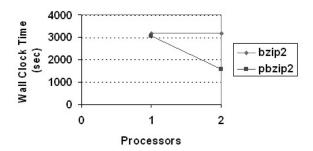


Figure 4. Wall clock time in seconds as a function of the number of processors to compress 1.83 GB on an Itanium2 900MHz system.

processors, 96.9% of linear speedup was achieved. A graph of the wall clock time as a function of the number of processors is shown in Figure 4.

5.3 SunFire 6800

The platform used in this experiment was a SunFire 6800 machine with twenty four 1.05 GHz UltraSPARC-III processors, 8MB L2 cache, and 96 GB of RAM. The software was compiled with gcc v3.2.3 and run on SunOS 5.9.

The pbzip2 program was run once with the -p1 switch for 1 processor, and tested again with the appropriate command line switches for 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, and 22 processors. The SunFire machine is shared amongst various researchers so it was not possible to get exclusive access to more than 22 processors.

The bzip2 program establishes a baseline for performance. Running pbzip2 on one processor showed that it was approximately the same speed as bzip2. It was observed that when running pbzip2 on multiple processors, near-optimal speedup was achieved. Using 18 processors, pbzip2 still performed at 98% of optimal speedup, and with 22 processors it was at 97.4% of optimal speedup. A graph of the wall clock time as a function of the number of processors is shown in Figure 5. The relative speedup as a function of the number of processors can be found in Figure 6.

6 Conclusion

In this paper, several methods of parallelizing the Burrows-Wheeler Transform are described. One method, the simultaneous processing of BWT blocks, was implemented and the performance of the algorithm was evaluated on several shared memory parallel architectures. The results showed that a significant, near-optimal speedup was achieved by using the pbzip2 algorithm on systems with multiple processors. This will greatly reduce the time it takes to com-

SunFire 6800 (22 x 1.05GHz)

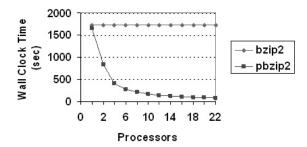


Figure 5. Wall clock time in seconds as a function of the number of processors to compress 1.83 GB on a SunFire 6800 system.

SunFire 6800 (22 x 1.05GHz)

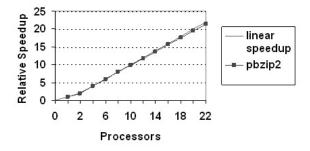


Figure 6. Relative speedup as a function of the number of processors to compress 1.83 GB on a SunFire 6800 system.

press large amounts of data while remaining fully compatible with the sequential version of bzip2.

Future work may consist of implementing pbzip2 on a distributed memory cluster architecture, and exploring the parallel sort method of parallelizing the BWT algorithm.

7 Acknowledgements

Special thanks to Dr. Frank Dehne and John Saal-waechter for their suggestions, and Ken Hines, HPCVL (www.hpcvl.org), and SHARCNET (www.sharcnet.com) for use of their computing resources.

References

- [1] M. Burrows and D. J. Wheeler. A block-sorting lossless data compression algorithm. Technical Report 124, 1994.
- [2] Brenton Chapin and Stephen R. Tate. Higher compression from the burrows-wheeler transform by modified sorting. In *Data Compression Conference*, page 532, 1998.
- [3] J. G. Cleary and I. H. Witten. Data compression using adaptive coding and partial string matching. *IEEE Transactions on Communications*, COM-32(4):396–402, April 1984.
- [4] Michelle Effros. Universal lossless source coding with the burrows wheeler transform. In *Data Compression Conference*, pages 178–187, 1999.
- [5] P. Fenwick. Block-sorting text compression final report, 1996.
- [6] P. Franaszek, J. Robinson, and J. Thomas. Parallel compression with cooperative dictionary construction. U.S. Patent No. 5,729,228, 1998.
- [7] R. Nigel Horspool. Improving LZW. In *Data Compression Conference*, pages 332–341, 1991.
- [8] Paul G. Howard and Jeffery Scott Vitter. Arithmetic coding for data compression. Technical Report Technical report DUKE-TR-1994-09, 1994.
- [9] Julian Seward. The bzip2 and libbzip2 official home page (http://sources.redhat.com/bzip2/), 2002.
- [10] Lynn M. Stauffer and Daniel S. Hirschberg. Parallel text compression - REVISED. Technical Report ICS-TR-91-44, 1993.
- [11] Lynn M. Stauffer and Daniel S. Hirschberg. Dictionary compression on the PRAM. Technical Report ICS-TR-94-07, 1994.
- [12] Jacob Ziv and Abraham Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3):337–343, 1977.