# The Engineering of a Compression Boosting Library:
# Theory vs Practice in BWT Compression

Paolo Ferragina[1,*], Raffaele Giancarlo[2,**], and Giovanni Manzini[3,***]

[1] Dipartimento di Informatica, Università di Pisa, Italy
`ferragina@di.unipi.it`
[2] Dipartimento di Matematica ed Applicazioni, Università di Palermo, Italy
`raffaele@math.unipa.it`
[3] Dipartimento di Informatica, Università del Piemonte Orientale, Italy
`manzini@mfn.unipmn.it`

**Abstract.** Data Compression is one of the most challenging arenas both for algorithm design and engineering. This is particularly true for Burrows and Wheeler Compression a technique that is important in itself and for the design of compressed indexes. There has been considerable debate on how to design and engineer compression algorithms based on the BWT paradigm. In particular, Move-to-Front Encoding is generally believed to be an "inefficient" part of the Burrows-Wheeler compression process. However, only recently two theoretically superior alternatives to Move-to-Front have been proposed, namely Compression Boosting and Wavelet Trees. The main contribution of this paper is to provide the first experimental comparison of these three techniques, giving a much needed methodological contribution to the current debate. We do so by providing a carefully engineered compression boosting library that can be used, on the one hand, to investigate the myriad new compression algorithms that can be based on boosting, and on the other hand, to make the first experimental assessment of how Move-to-Front behaves with respect to its recently proposed competitors. The main conclusion is that Boosting, Wavelet Trees and Move-to-Front yield quite close compression performance. Finally, our extensive experimental study of boosting technique brings to light a new fact overlooked in 10 years of experiments in the area: a fast adapting order-zero compressor is enough to provide state of the art BWT compression by simply compressing the run length encoded transform. In other words, Move-to-Front, Wavelet Trees, and Boosters can all be by-passed by a fast learner.

## 1   Introduction

In the quest for the ultimate data compressor, Algorithmic Theory and Engineering go hand in hand. This point is well illustrated by the amount of results and implementations originated by the fundamental results by Lempel and Ziv. A more recent example is provided by the fundamental contributions given by Burrows and Wheeler to data compression [3], via their transform (denoted for short bwt). In their seminal paper Burrows and Wheeler proposed to compress the output of the bwt using Move-to-Front Encoding (shortly mtf), followed by an order zero compressor (usually Arithmetic or Huffman coding). As pointed out by Fenwick [5] in the first systematic study of that new type of compression, the technique is so powerful that it yields nearly state-of-the-art compression results without any particularly sophisticated engineering of the coding step. This should be contrasted with PPM-based compressors that involve quite a bit of engineering. From that point on, the research on bwt compression has focused on two aspects: faster bwt computation, and the identification and exploitation of potential inefficiencies in the use of mtf. While substantial progress has been made on the first point, both theoretically and experimentally (e.g. [2, 17]), the second point experienced a plethora of heuristically-designed proposals (see [1, 4] and references therein) which improved over the original proposal but often lacked of analytical justification.

Recently, two theoretical results [7, 8] have shed new light on the role of mtf within the bwt-based compression paradigm, paving the way to the (*analytically justified*) design of more powerful bwt-based compressors. In particular, [8] proposed a new technique, named *compression boosting*, that fully uses the power of bwt to show that the performance of *any* order zero compressor can be automatically, and optimally, boosted to higher order entropy compression. On the other hand, [7] proved that combining the bwt with the Wavelet Tree data structure [10] we can achieve high-order entropy bounds without using mtf or the boosting technique. At the same time, a novel and very recent analysis of classic bwt compression [12] showed that mtf may not be as inefficient as initially thought. Summing this with the fact that the theoretical results in [7, 8] require some sophisticated algorithmic machinery, it is not at all clear how much computational/compression *gain* can be achieved by shaving off the mtf-step from the bwt-based compressors.

The above is the main question addressed in the present paper, whose key contribution is first of all *methodological*. We provide the first carefully engineered compression boosting library that can be used, on the one hand, to investigate the myriad new compression algorithms that can be based on boosting, and on the other hand, to make the first experimental assessment of how mtf behaves with respect to its recently proposed competitors: Boosting and Wavelet Trees. The boosting library is available under the GPL license at the page http://www.mfn.unipmn.it/~manzini/boosting and it is highly modular in the sense that it can be used to create a powerful high order compressor even without any knowledge of the bwt.

In order to highlight our additional technical contributions, we need to recall a few facts about compression boosting [8]. Additional details are given in Section 3. The boosting technique builds upon three main ingredients: bwt, the Suffix Tree data structure, and a greedy algorithm to process them. Specifically, it is shown that there exists a proper partition of the bwt of a string $s$ exhibiting a deep combinatorial relation with the $k$-th order entropy of $s$. That partition can be identified via a greedy processing of the suffix tree of $s$. The final compressed string is then obtained by compressing individually each substring of the partition by means of the base (order zero) compressor A we wish to boost. The proper design of a compression booster is a bit trickier than it sounds:

(**A**) The greedy algorithm alluded to before is a bottom up visit of the suffix tree. In practice, on large files, the memory requirements for the construction of the suffix tree would be prohibitively large. We use suffix arrays instead and procedures that efficiently simulate the bottom up visit of the suffix tree [13].

(**B**) Given the algorithm A we wish to boost, we also need an objective function that estimates how well A compresses a given string. In [8], the objective function is given in terms of two parameters $\lambda$ and $\mu$, and the order zero empirical entropy of the string (see Section 3 for details). In practice, $\lambda$ and $\mu$ may either be not available or be too conservative. This point is discussed in Section 4, where we propose two cost models and the relative objective functions.

(**C**) Another important aspect of the boosting process is the ability of the algorithm A to quickly adapt to the statistics of a string to be compressed. Faster adaptation means better compression. This learning process is usually governed by parameters establishing how fast A "forgets the past". We limit our experimentation to range coding and arithmetic coding. The somewhat intuitive, yet surprising, results are reported in Section 5 and outlined in point (**F**) below.

Using our library we have compared the performance of the compression booster against bwt compressors based on mtf (e.g. Bzip2 [19] and variants), bwt compressors based on Wavelet Trees (e.g., Wzip [9]), and state-of-the-art PPM compressors (e.g. PPMd [21]). We show that:

(**D**) As predicted by Theory [8], boosting is superior to classic bwt approaches that use mtf in terms of compression ratio but not by much. It is also slower, as it is to be expected, because of the significant time cost for building the optimal bwt-partition (as observed in **B**). Therefore, those results give a strong indication that mtf may actually be a time-efficient way to effectively "approximate" the optimal partition computed by the boosting technique.

(**E**) As predicted by Theory [7, 10, 11], the simple combination of bwt with Wavelet Trees is effective both in time and compression ratio, and does not benefit from the use of the booster. However, the Wavelet Tree approach is outperformed by classic bwt approaches that use mtf. This further confirms the effectiveness in time and compression ratio of mtf, and leaves open the problem of investigating the more powerful approach proposed in [7], namely *Generalized Wavelet Trees*, which are based on sophisticated combinations of binary (like,

Run Length encoders) versus non-binary (like, Huffman or Arithmetic encoders) compressors and Wavelet Trees of properly-designed shapes.

(**F**) The experiments performed to estimate the best adaptation parameters for range and arithmetic coding show clearly that a fast adaptation yields state-of-the-art compression by simply compressing a run length encoded bwt. This is somewhat intuitive, yet surprising: to our knowledge no one observed experimentally the superiority of this strategy w.r.t. mtf, and no theoretical analysis has explained or suggested such behavior. Moreover, this result comes from the stronger finding that for a fast adapting range coder the optimal partition coming out of the booster is the bwt itself (data not shown, due to space limitations). That is, the strategy is optimal with respect to the boosting paradigm.

(**G**) All the bwt-based compressors we tested were inferior, in terms of compression ratio, to the highly engineered PPMd tool. The principle behind bwt and PPM techniques is the same: discover and encode according to the "best" contexts. However, bwt-based algorithms have the advantage of knowing the entire string, while PPMd "discovers" good contexts on-line. Yet bwt-based algorithms do not perform as well. This yields an extremely intriguing engineering problem for data compression practitioners. Note that there is a very good reason to stick with bwt-based compressors instead of embracing the, apparently superior, PPM-based compressors: the reason is that bwt-based compressors are a key tool for the construction of compressed indices which (informally) are compressed files offering the additional capability of very fast full-text search (see [18] for formal definitions and a comprehensive survey).

In conclusion our experiments show that Boosting, Wavelet Trees and mtf yield quite close compression performance. However, the boosting technique appears to be more robust and works well even with less effective order zero compressors (such as Huffman coding). Moreover, when used with range/arithmetic coding the boosting technique yields excellent compression somewhat irrespective of how fast the order-zero compressor adapts to the statistics of the string. These positive features are achieved using more resources (time and space) during compression: nevertheless our results show that a careful implementation of boosting can handle efficiently even very large files.

## 2   Background and Notation

Let $s$ be a string over the alphabet $\Sigma = \{a_1, \ldots, a_h\}$ and, for each $a_i \in \Sigma$, let $n_i$ be the number of occurrences of $a_i$ in $s$. The 0-*th order empirical entropy* of the string $s$ is defined as[1] $H_0(s) = -\sum_{i=1}^{h}(n_i/|s|)\log(n_i/|s|)$. It is well known that $H_0$ is the maximum compression we can achieve using a fixed codeword for each alphabet symbol. We can achieve a greater compression if the codeword we use for each symbol depends on the $k$ symbols preceding it, since the maximum compression is now bounded by the $k$-th order entropy $H_k(s)$ (see [15] for the formal definition). For highly compressible strings, $|s|\,H_k(s)$ fails to provide a

---

[1] We assume that all logarithms are taken to the base 2 and $0\log 0 = 0$.

```
$ mississipp i
i $mississip p
i ppi$missis s
i ssippi$mis s
i ssissippi$ m
m ississippi $
p i$mississi p
p pi$mississ i
s ippi$missi s
s issippi$mi s
s sippi$miss i
s sissippi$m i
```
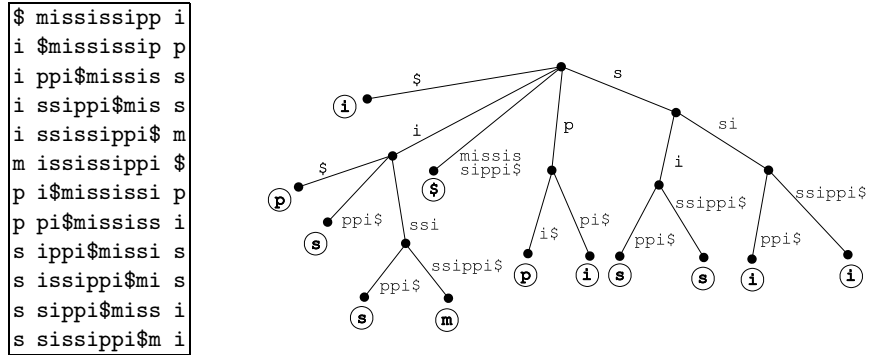


**Fig. 1.** The bwt matrix (left) and the suffix tree (right) for the string $s = $ mississippi$. Note that the output of the bwt is the last column of the bwt matrix, i.e., ipssm$pissii.

reasonable bound to the performance of compression algorithms (see discussion in [8,15]). For that reason, [15] introduced the notion of 0-*th order modified empirical* entropy $H_0^*(s)$ which has the property that if $|s| > 0$, $|s|H_0^*(s)$ is at least equal to the number of bits needed to write down the length of $s$ in binary. The *k-th order modified empirical entropy* $H_k^*$ is then defined in terms of $H_0^*$ as the maximum compression we can achieve by looking at *no more than k* symbols preceding the one to be compressed.

Given a string $s$, the Burrows-Wheeler transform (bwt for short) consists of three basic steps: (1) append to the end of $s$ a special symbol $ smaller than any other symbol in $\Sigma$; (2) form a *conceptual* matrix $\mathcal{M}$ whose rows are the cyclic shifts of the string $s$, sorted in lexicographic order; (3) construct the transformed text $\hat{s} = \text{bwt}(s)$ by taking the last column of $\mathcal{M}$ (see Fig. 1). Although it is not obvious, from $\hat{s}$ we can always recover $s$, see [3] for details. The power of the bwt rests on the fact that equal contexts (substrings) of $s$ are grouped together resulting in a few clusters of distinct symbols in $\text{bwt}(s)$. That clustering makes $\text{bwt}(s)$ a better string to compress than $s$. In their seminal paper Burrows and Wheeler proposed to compress the output of the bwt using Move-to-Front Encoding[2] (shortly mtf), followed by an order zero compressor (Arithmetic or Huffman coding). In [12] it is shown that if we use an order zero compressor A such that for any string $x$ we have $|A(x)| \leq |x|H_0(x) + c|x|$, then the bwt followed by mtf, followed by A produces an output bounded by

$$\mu|s|H_k(s) + (\log \zeta(\mu) + c)|s| + \log |s| + \mu g_k \qquad (1)$$

where $\zeta$ is the Riemann zeta function. The bound (1) holds for any $k \geq 0$ and $\mu > 1$. In [15] it is shown that if we use Run Length Encoding (shortly rle) between mtf and the order zero compressor, the output is bounded by

$$(5 + \epsilon)|s|H_k^*(s) + \log_2 |s| + g_k' \qquad (2)$$

---

[2] Move-to-Front transforms the input encoding each symbol with the number of distinct symbols seen since its last occurrence, see [3] for details.

for any $k \geq 0$ and $\epsilon \approx 10^{-2}$. The bottom line is that combining the Burrows-Wheeler transform with $\mathtt{mtf}$ and an order zero compressor we can achieve the $k$-th order entropy, $H_k$ or $H_k^*$, simultaneously for any $k \geq 0$. Note however, that the coefficient in front of the $k$-th order entropy in (1) and (2) is greater than 1 whereas we are assuming that $\mathtt{A}$ achieves $H_0$ without any multiplicative constant. This means that there is a small inefficiency as we go from $H_0$ and $H_0^*$ to $H_k$ and $H_k^*$. It is an open question whether this inefficiency can be removed with a more detailed analysis or is inherent in the use of Move-to-Front encoding.

## 3   A BWT-Based Compression Booster

Recently [8] has described a $\mathtt{bwt}$-based compression booster that, starting from an order zero compressor, achieves the $k$-th order entropy without the inefficiency found in the $\mathtt{mtf}$-based approach. In this section we quickly review how the boosting algorithm works; the details and proofs can be found in [8].

A crucial ingredient of the compression booster is the relationship between the $\mathtt{bwt}$ matrix and the suffix tree data structure. Let $\mathcal{T}$ denote the suffix tree of the string $s\$$. $\mathcal{T}$ has $|s| + 1$ leaves, one per suffix of $s\$$, and edges labeled with substrings of $s\$$ (see Figure 1). Any node $u$ of $\mathcal{T}$ has *implicitly associated* a substring of $s\$$, given by the concatenation of the edge labels on the downward path from the root of $\mathcal{T}$ to $u$. In that implicit association, the leaves of $\mathcal{T}$ correspond to the suffixes of $s\$$. We assume that the suffix tree edges are sorted lexicographically. Since each row of the $\mathtt{bwt}$ matrix is prefixed by one suffix of $s\$$ and rows are lexicographically sorted, the $i$-th leaf (counting from the left) of the suffix tree corresponds to the $i$-th row of the $\mathtt{bwt}$ matrix. We associate the $i$-th leaf of $\mathcal{T}$ with the $i$-th symbol of the string $\hat{s} = \mathtt{bwt}(s)$. The symbol associated to the leaf $v$ is thus the symbol preceding in $s$ the substring of $s\$$ associated with $v$. Such symbols are represented inside circles in Fig. 1. If we write $\hat{\ell}_i$ to denote the symbol associated to the $i$-th leaf, from the above discussion, it follows that $\hat{s} = \hat{\ell}_1 \hat{\ell}_2 \cdots \hat{\ell}_{|s|+1}$ (see Fig. 1 for an example).

For any suffix tree node $u$, let $\hat{s}\langle u \rangle$ denote the substring of $\hat{s}$ obtained concatenating, from left to right, the symbols associated to the leaves descending from node $u$. We say that a subset $\mathcal{L}$ of $\mathcal{T}$'s nodes is a *leaf cover* if every leaf of the suffix tree has a *unique* ancestor in $\mathcal{L}$. Any leaf cover $\mathcal{L} = \{u_1, \ldots, u_p\}$ naturally induces a partition of the leaves of $\mathcal{T}$ namely $\hat{s}\langle u_1 \rangle, \ldots, \hat{s}\langle u_p \rangle$. Because of the relationship between $\mathcal{T}$ and the $\mathtt{bwt}$ matrix this is also a partition of $\hat{s}$.

Let $C$ denote a function which associates to every string $x$ over $\Sigma \cup \{\$\}$ the positive real value $C(x)$. For any leaf cover $\mathcal{L}$, we define its cost as: $C(\mathcal{L}) = \sum_{u \in \mathcal{L}} C(\hat{s}\langle u \rangle)$. In [8] it is shown a linear time greedy algorithm that computes a leaf cover $\mathcal{L}_{\min}$ of minimum cost. That is, $\mathcal{L}_{\min}$ is such that $C(\mathcal{L}_{\min}) \leq C(\mathcal{L})$, for any leaf cover $\mathcal{L}$. $\mathcal{L}_{\min}$ is called an *optimal leaf cover* and we say that $\mathcal{L}_{\min}$ induces an *optimal partition* of $\hat{s}$ with respect to the cost function $C$. The relevance of $\mathcal{L}_{\min}$ for achieving the $k$-th order entropy derives by the following Theorem [8].

**Theorem 1.** *Let* $\mathtt{A}$ *denote an order zero compressor such that for any string* $x$ $|\mathtt{A}(x)| \leq \lambda |x| \, H_0^*(x) + \mu$ *where* $\lambda$ *and* $\mu$ *are constants. Let* $\mathcal{L}_{\min}$ *denote an*

*optimal partition of $\hat{s}$ with respect to $C(x) = \lambda|x| H_0^*(x) + \mu$. If we use algorithm* A *to compress the substrings of the optimal partition induced by $\mathcal{L}_{\min}$, the overall output size is bounded by $\lambda|s| H_k^*(s) + g_k$ bits for any $k \geq 0$, where $g_k$ only depends on the alphabet size $|\Sigma|$. A similar result holds for $H_k$ as well.*    □

## 4  The Compression Boosting Library

The efficient implementation of the compression booster algorithm is a non trivial engineering task. The main challenge is avoiding the explicit construction of the suffix tree which would require an unpractically large amount of working memory. We now detail our implementation discussing its space requirements in the "real world" model where we assume that every character takes one byte and every integer takes 4 bytes. Let $n = |s|$. We first compute the suffix array of $s$ using the ds algorithm [17] that has a peak memory usage of only $5.03n$ bytes: $n$ bytes for the text, $4n$ for the suffix array, and $0.03n$ working space.

Given the suffix array we compute and store $\hat{s} = \mathtt{bwt}(s)$ using $n$ bytes. The greedy algorithm computing the optimal partition of $\hat{s}$ consists of a properly defined post-order visit of the suffix tree of $s$. To avoid the explicit construction of the suffix tree we use the technique from [13] that allows one to emulate the post-order visit of the suffix tree using the Longest Common Prefix (shortly LCP) array. Thus, we use the Lcp6 algorithm from [16] for computing in $O(n)$ time the LCP array given $s$, $\hat{s}$, and the suffix array. This algorithm overwrites the LCP array over the suffix array and has a peak space usage of $(6 + \delta)n$ bytes. The parameter $\delta$ is at most 4 and is bounded also by $|\Sigma|^k/n + 2H_k(s)$ for any $k \geq 0$. This means that the space usage is smaller for highly compressible inputs.

Having computed the LCP array we can discard the input string $s$; thus at this stage we are only storing $\hat{s}$ and the LCP array for a total space usage of $5n$ bytes. The computation of the optimal partition using the technique in [13] reduces to a left to right scan of the LCP array. This allows us to store the endpoints of intervals of the optimal partition in the same memory used for the LCP array (that is, overwriting the LCP array). Thus the only additional memory used during the "emulated" suffix tree visit is the space used to store the stack of the suffix tree nodes whose visit has started but not yet finished. This space could be $\Theta(n)$ in the worst case, but in practice is much smaller than $n$ bytes overall.

**Cost models.** An important issue in the implementation of the compression booster is the choice of the parameters $\lambda$ and $\mu$ in the cost function $C(x) = \lambda|x| H_0^*(x) + \mu$ of Theorem 1. Given a compressor A, theory dictates that $\lambda$ and $\mu$ be chosen so that $|\mathtt{A}(x)| \leq C(x)$ for any string $x$. However, if we strictly enforce this condition it is possible that for many strings $x$ we have $|\mathtt{A}(x)| \ll C(x)$. Since the optimal partitioning is computed minimizing $C(\mathcal{L}_{\min})$, if $C(x)$ is "too far" from $|\mathtt{A}(x)|$ we could end up with a partition which does not exploit the full potential of the compressor A. To evaluate this phenomenon our boosting library supports two different cost models. In addition to the "entropy bound" model

outlined above, we provide a "real cost" model in which the optimal partition is computed with respect to the cost $C(x) = |\mathtt{A}(x)|$. Using the "real cost" model we get the best possible compression that we can achieve using the compressor $\mathtt{A}$. The drawback of this model is that the computation of the optimal partition no longer takes linear time. The time cost might be quadratic in the worst case, although the experimental results show that the overall running time usually increases only by a factor 1.5.

**User interface.** Our library provides a simple interface to boost the performance of an arbitrary compressor using either mtf or the optimal partitioning strategy outlined in Sect. 3. This can be done even without any knowledge of the Burrows-Wheeler transform! The user simply needs to provide compression and decompression procedures and, for the computation of the optimal partition, a procedure evaluating the cost function $C(x)$ (see [6] for details).

## 5 Experimental Results

Using the boosting library described in the previous section we have implemented several bwt-based compressors. By means of extensive experiments we tried to assess to what extent mtf and the boosting algorithm are able to turn a generic order zero compressor into a state of the art compressor. We ran all experiments on a 2.6 GHz Pentium 4 CPU with 1.5 GB of main memory running Fedora Linux. All code was written in C and compiled using gcc Ver. 3.2.2. As a testbed we used the collection of files introduced in [17] for testing suffix array construction algorithms.

The following are the algorithms tested in our experiments.

Bzip2 is the well known tool based on the bwt developed by Julian Seward [19]. Bzip2 splits the input file into blocks of size 900Kb and computes the bwt followed by mtf on each block. The actual compression is done using rle0[3] followed by Multiple-Table Huffman coding [22].

MtfRleMth executes the same steps as Bzip2 operating on the whole input instead that on fixed length blocks.

MtfRleRc. The earliest versions of Bzip2 used arithmetic coding instead of multiple-table Huffman. Recently, range coding has been (re)discovered as a patent-free alternative to arithmetic coding. Range coding and arithmetic coding are based on similar concepts and achieve similar compression. MtfRleRc compresses the bwt using mtf followed by rle0, followed by range coding (we used the code from [14]). Note that MtfRleRc is identical to MtfRleMth except that, instead of Multiple-table Huffman coding, it uses range coding.

RleRc compresses the bwt using rle followed by range coding.

---

[3] We use rle to denote the run length encoding of the runs of any character, while we use rle0 to denote the run length encoding only of the runs of zeros. If a string was produced by mtf, rle0 is the natural choice because of the massive presence of 0-runs as observed by Fenwick [5].

BoostRleRc is the boosting algorithm applied to the compressor consisting of rle followed by range coding. Note that the difference in compression between RleRc and BoostRleRc gives the "added value" of the use of the booster.

MtfRleAc, RleAc, BoostRleAc are analogous respectively to MtfRleRc, RleRc, BoostRleRc except that they use the arithmetic coding routines from [23] instead of range coding.

MtfRleHuff, RleHuff, BoostRleHuff are analogous respectively to MtfRleRc, RleRc, BoostRleRc except that they use Huffman coding instead of range coding. Note that MtfRleHuff differs from MtfRleMth in that the former uses a single Huffman table whereas the latter uses up to six tables for the same file.

Wavelet. This algorithm computes the bwt of the whole input and compresses the resulting string using a wavelet tree [10]. The importance of wavelet trees stems from the fact that they have been used for the design of efficient bwt-based compressed indices [18] and that they also achieve the $k$-th order entropy for any $k \geq 0$. More precisely, from [7] follows that for a string $s$ over the alphabet $\Sigma$ the output size of Wavelet is bounded by $4|s| H_k^*(s) + 6|\Sigma|^{k+1} \log(|s|)$ bits.

BoostWav is an implementation of the boosting algorithm applied to the wavelet tree encoder using the "real cost" model. Thus the difference between Wavelet and BoostWav is that the former builds one wavelet tree on the whole bwt, whereas the latter finds an optimal partition of the bwt and builds one wavelet tree on each substring of the optimal partition. Again, the difference in compression between Wavelet and BoostWav is the "added value" of the booster.

PPMd is an implementation of the ppm encoder by Dmitry Shkarin [21, 20] which is the current state of the art for PPM compression. In our tests we used PPMd at its maximum strength, that is using a model of order 16 and 256Mb of working memory.

**Range/arithmetic coding variants.** The behavior of range and arithmetic coding depends on two parameters: MaxFreq and Increment. The ratio between these two values essentially controls how quickly the coding "adapts" to the new statistics. For range coding we set MaxFreq = 65536 (the largest possible value) and we experimented with three different values of Increment. Setting Increment = 256 we get a range coder with FAST adaptation, with Increment = 32 we get a range coder with MEDIUM adaptation, and finally setting Increment = 8 we get a range coder with SLOW adaptation. For arithmetic coding we set MaxFreq = 16383 (the largest possible value) and Increment = 64 obtaining therefore a FAST adaptation.

**Compression ratio.** Figure 2 reports the average compression ratio (in bits per symbol) and average (de)compression time (microseconds per symbol) for all the algorithms mentioned above. Looking at the average compression ratio we can see that both mtf and the boosting algorithm do a good job in transforming an order zero compressor into a state-of-the-art compressor. However, our data show some unexpected behaviors. Considering the three version of range coding (with FAST, MEDIUM, and SLOW adaptation) we see that mtf achieves the best compression using MEDIUM adaptation whereas the boosting algorithm "prefers"

| | averg | ctime | dtime |
|---|---|---|---|
| Bzip2 | 1.424 | 0.53 | 0.14 |
| MtfRleMth | 1.167 | 0.96 | 0.46 |
| RleAc FAST | 1.126 | 0.97 | 0.59 |
| MtfRleAc FAST | 1.158 | 0.94 | 0.53 |
| BoostRleAc FAST RC | 1.125 | 7.43 | 0.59 |
| RleHuff | 1.596 | 0.89 | 0.47 |
| MtfRleHuff | 1.230 | 0.95 | 0.46 |
| BoostRleHuff RC | 1.195 | 5.04 | 0.45 |
| BoostRleHuff EB | 1.229 | 2.96 | 0.45 |
| Wavelet | 1.230 | 0.96 | 1.01 |
| BoostWav RC | 1.229 | 3.55 | 0.94 |
| PPMd | 1.080 | 0.60 | 0.66 |

| | averg | ctime | dtime |
|---|---|---|---|
| RleRc FAST | 1.129 | 0.90 | 0.48 |
| MtfRleRc FAST | 1.161 | 0.90 | 0.48 |
| BoostRleRc FAST RC | 1.129 | 4.11 | 0.48 |
| BoostRleRc FAST EB | 1.134 | 3.04 | 0.48 |
| RleRc MED. | 1.171 | 0.89 | 0.48 |
| MtfRleRc MED. | 1.153 | 0.96 | 0.48 |
| BoostRleRc MED. RC | 1.152 | 4.13 | 0.49 |
| BoostRleRc MED. EB | 1.158 | 3.02 | 0.48 |
| RleRc SLOW | 1.245 | 0.90 | 0.48 |
| MtfRleRc SLOW | 1.164 | 0.90 | 0.48 |
| BoostRleRc SLOW RC | 1.175 | 4.12 | 0.48 |
| BoostRleRc SLOW EB | 1.194 | 3.02 | 0.48 |

**Fig. 2.** Experimental results for the collection of files introduced in [17]. For each algorithm we report the average compression in bits per symbol and the average compression and decompression time in microseconds per symbol. The RC and EB acronyms indicate the cost model ("real cost" or "entropy bound") used by the booster.

| | running time | | | | | peak memory | |
|---|---|---|---|---|---|---|---|
| | bwt | lcp | visit | cmpr | total | lcp | visit |
| *sprot* | 0.70 | 0.60 | 1.67 | 0.11 | 3.08 | 7.01 | 5.00 |
| *rfc* | 0.60 | 0.51 | 2.35 | 0.11 | 3.57 | 6.86 | 5.00 |
| *howto* | 0.50 | 0.45 | 2.83 | 0.15 | 3.93 | 7.29 | 5.01 |
| *reut* | 1.24 | 0.55 | 1.92 | 0.08 | 3.79 | 6.58 | 5.00 |
| *linux* | 0.52 | 0.42 | 3.39 | 0.12 | 4.46 | 6.88 | 5.04 |
| *jdk13* | 1.15 | 0.40 | 2.10 | 0.05 | 3.70 | 6.26 | 5.00 |
| *etext* | 0.75 | 0.63 | 2.65 | 0.16 | 4.19 | 7.57 | 5.00 |
| *chr22* | 0.49 | 0.54 | 6.33 | 0.17 | 7.53 | 8.34 | 5.49 |
| *gcc* | 0.85 | 0.40 | 3.00 | 0.10 | 4.36 | 6.75 | 5.07 |
| *w3c* | 1.10 | 0.43 | 3.18 | 0.06 | 4.78 | 6.31 | 5.01 |

| | running time | | | | |
|---|---|---|---|---|---|
| | bwt | lcp | visit | cmpr | total |
| *sprot* | 0.70 | 0.59 | 1.23 | 0.11 | 2.63 |
| *rfc* | 0.60 | 0.51 | 1.52 | 0.11 | 2.74 |
| *howto* | 0.50 | 0.46 | 1.86 | 0.15 | 2.96 |
| *reut* | 1.24 | 0.56 | 1.32 | 0.08 | 3.19 |
| *linux* | 0.52 | 0.42 | 2.17 | 0.12 | 3.23 |
| *jdk13* | 1.15 | 0.40 | 1.48 | 0.05 | 3.08 |
| *etext* | 0.75 | 0.64 | 1.59 | 0.16 | 3.14 |
| *chr22* | 0.49 | 0.54 | 0.89 | 0.18 | 2.10 |
| *gcc* | 0.86 | 0.40 | 1.64 | 0.10 | 3.00 |
| *w3c* | 1.10 | 0.43 | 2.34 | 0.06 | 3.94 |

**Fig. 3.** Running time and peak memory usage for the various stages of the BoostRleRc (MEDIUM adaptation) algorithm using the "real cost" model (left) and the "entropy bound" model (right, the table only shows running times since the memory usage is the same as for the "real cost" model). The running times of the four basic steps (bwt computation, LCP array computation, optimal partition computation via suffix tree visit, actual compression using range coding) and the total running time are given in microseconds per input byte. The peak memory usage is given for the LCP array computation and the suffix tree visit which are the steps using more memory. Memory usage is reported as number of used bytes per input byte.

FAST adaptation. It is also remarkable that RleRc with FAST adaptation achieves a very good compression, better indeed that mtf combined with any version of range coding (and the same is true for RleAc FAST). This means that the bwt can be compressed efficiently using rle and an order zero encoder that quickly

adapts to the new statistics. This is somewhat intuitive, but to our knowledge no one observed experimentally the superiority of this strategy w.r.t. mtf, and no theoretical analysis has explained or suggested such behavior. Overall the data show that the boosting algorithm is superior to mtf in terms of compression ratio and it is also more robust in the sense that it works well even with less effective order zero compressors (for example Huffman coding). This superiority is however paid in terms of running time as discussed below.

**Running time.** The data in Figure 2 show that for range coding the boosting algorithm with the "real cost" model is between 4 and 5 times slower than mtf in compression while there is no significant difference in decompression. For arithmetic and Huffman coding the ratio is even higher. Using the "entropy bound" model the compression time decreases significantly and there is a corresponding loss in compression efficiency. Summing up, mtf and the boosting algorithm (with the two different cost models) offer three different trade offs between compression ratio and compression time: the user can choose the one most suitable for the application at hand. Figure 3 reports the resource usage of the various stages of the boosting algorithm. We can see that the most time consuming step is the optimal partition computation via the suffix tree visit both in the "real cost" and "entropy bound" models. Note also that the peak memory usage is achieved during the LCP array computation.

**Wavelet tree performance.** The data in Figure 2 show that the algorithms Wavelet and BoostWav roughly achieve the same compression as the algorithms based on Huffman coding (RleHuff and BoostRleHuff) and are inferior to the algorithms based on range/arithmetic encoding. We point out that the similar compression ratio of Wavelet and BoostWav provide an experimental validation of the theoretical analysis of [7] which states that even using a single wavelet tree—as in the algorithm Wavelet—we already achieve the $k$-th order entropy.

**PPMd performance.** The results in Figure 2 show that PPMd outperforms all other compressors. Additional tests on the files of the Canterbury corpus (see [6]) show that the Weighted Frequency Count algorithm from [1] (which is based on the bwt) compresses better than mtf, boosting, and wavelet tree algorithms. This suggests that in the field of (bwt) compression Theory is currently a step behind Practice. Although we emphasize that for the construction of compressed indexes it is essential to have simple and efficient bwt-based algorithms whose performance are theoretically guaranteed, we take these results as a stimulus for further research!

# References

1. J. Abel. Post BWT stages of the Burrows-Wheeler compression algorithm. Submitted. See also "A fast and efficient post BWT-stage for the Burrows-Wheeler compression algorithm". Proc. *IEEE DCC*, 2005, pag. 449.
2. S. Burkhardt and J. Kärkkäinen. Fast lightweight suffix array construction and checking. In *Proc. 14th Symposium on Combinatorial Pattern Matching (CPM '03)*, pages 55–69. Springer-Verlag LNCS n. 2676, 2003.

3. M. Burrows and D. Wheeler. A block sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, 1994.

4. S. Deorowicz. Context exhumation after the BurrowsWheeler transform. *Information Processing Letters*, 95:313–320, 2005.

5. P. Fenwick. Block sorting text compression — final report. Technical Report 130, Dept. of Computer Science, The University of Auckland New Zeland, 1996.

6. P. Ferragina, R. Giancarlo, and G. Manzini. The engineering of a compression boosting library: Theory vs practice in BWT compression. Technical Report TR-INF-2006-06-03-UNIPMN, http://www.di.unipmn.it, 2006.

7. P. Ferragina, R. Giancarlo, and G. Manzini. The myriad virtues of wavelet trees. In *Proc. of International Colloquium on Automata and Languages (ICALP)*, pages 561–572. Springer Verlag LNCS n. 4051, 2006.

8. P. Ferragina, R. Giancarlo, G. Manzini, and M. Sciortino. Boosting textual compression in optimal linear time. *Journal of the ACM*, 52:688–713, 2005.

9. L. Foschini, R. Grossi, A. Gupta, and J. Vitter. Fast compression with a static model in high order entropy. In *IEEE DCC*, pages 62–71. IEEE Computer Society TCC, 2004.

10. R. Grossi, A. Gupta, and J. Vitter. High-order entropy-compressed text indexes. In *Proc. 14th Annual ACM-SIAM Symp. on Discrete Algorithms (SODA '03)*, pages 841–850, 2003.

11. R. Grossi, A. Gupta, and J. Vitter. When indexing equals compression: Experiments on compressing suffix arrays and applications. In *Proc. 15th Annual ACM-SIAM Symp. on Discrete Algorithms (SODA '04)*, pages 636–645, 2004.

12. H. Kaplan, S. Landau, and E. Verbin. A simpler analysis of Burrows-Wheeler based compression. In *Proc. of the 17th Symposium on Combinatorial Pattern Matching (CPM '06)*. Springer-Verlag LNCS, 2006.

13. T. Kasai, G. Lee, H. Arimura, S. Arikawa, and K. Park. Linear-time longest-common-prefix computation in suffix arrays and its applications. In *Proc. 12th Symposium on Combinatorial Pattern Matching (CPM '01)*, pages 181–192. Springer-Verlag LNCS n. 2089, 2001.

14. M. Lundqvist. Carryless range coding. http://hem.spray.se/mikael.lundqvist/.

15. G. Manzini. An analysis of the Burrows-Wheeler transform. *Journal of the ACM*, 48(3):407–430, 2001.

16. G. Manzini. Two space saving tricks for linear time LCP computation. In *Proc. of 9th Scandinavian Workshop on Algorithm Theory (SWAT '04)*, pages 372–383. Springer-Verlag LNCS n. 3111, 2004.

17. G. Manzini and P. Ferragina. Engineering a lightweight suffix array construction algorithm. *Algorithmica*, 40:33–50, 2004.

18. G. Navarro and V. Mäkinen. Compressed full text indexes. Technical Report TR/DCC-2006-6, Dept. of Computer Science, University of Chile, 2006.

19. J. Seward. The BZIP2 home page, 2006. http://www.bzip.org.

20. D. Shkarin. PPMd compressor Ver. J. http://www.compression.ru/ds/.

21. D. Shkarin. PPM: One step to practicality. In *IEEE Data Compression Conference*, pages 202–211, 2002.

22. D. Wheeler. Improving Huffman coding, 1997. ftp://ftp.cl.cam.ac.uk/users/djw3/huff.ps.

23. I. H. Witten, R. M. Neal, and J. G. Cleary. Arithmetic coding for data compression. *Communications of the ACM*, 30(6):520–540, 1987.