

# Improving Run Length Encoding (RLE) on bit level by preprocessing

## Bachelorarbeit

zur Erlangung des akademischen Grades  
Bachelor of Science (B.Sc.)

Universität Trier  
FB IV - Informatikwissenschaften  
Lehrstuhl für Informatik I

Gutachter:	Prof. Dr. Ingo J. Timm xxxxxxxxxx
Betreuer:	xxxxxxxxxx

Vorgelegt am xx.xx.xxxx von:

Sven Fiergolla  
Am Deimelberg 30  
54295 Trier  
sven.fiergolla@gmail.com  
Matr.-Nr. 1252732

# Abstract

Hier steht eine Kurzzusammenfassung (Abstract) der Arbeit. Stellen Sie kurz und präzise Ziel und Gegenstand der Arbeit, die angewendeten Methoden, sowie die Ergebnisse der Arbeit dar. Halten Sie dabei die ersten Punkten eher kurz und fokussieren Sie die Ergebnisse. Bewerten Sie auch die Ergebnissen und ordnen Sie diese in den Kontext ein.

Die Kurzzusammenfassung sollte maximal 1 Seite lang sein.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Problem statement . . . . .	1
1.3	Main Objective . . . . .	2
1.4	Structure of this work . . . . .	2
<b>2</b>	<b>Principles of compression</b>	<b>3</b>
2.1	Compression and Encoding fundamentals . . . . .	3
2.1.1	Information Theory and Entropy . . . . .	3
2.1.2	General Analysis . . . . .	4
2.1.3	Probabilistic Coding . . . . .	4
2.1.4	Dictionary Coding . . . . .	5
2.1.5	Irreversible Compression . . . . .	5
2.2	Run Length Coding . . . . .	5
2.2.1	The history . . . . .	6
2.2.2	Limitations . . . . .	6
2.2.3	Run Length Encoding today . . . . .	7
2.3	Prefix Coding . . . . .	7
2.3.1	Huffman Coding . . . . .	7
2.4	Other methods . . . . .	7
2.4.1	Burrows-Wheeler-Transformation . . . . .	8
2.4.2	Inverse Burrows-Wheeler-Transformation . . . . .	8
2.5	State of the Art . . . . .	8
2.6	Limits of compression . . . . .	8
<b>3</b>	<b>Analysis</b>	<b>10</b>
3.1	Initial Findings . . . . .	10
3.2	Possible Improvements by Preprocessing . . . . .	12
3.2.1	Burrows-Wheeler-Transformation . . . . .	12
3.2.2	Vertical byte reading . . . . .	12
3.2.3	Byte remapping . . . . .	13
3.2.4	Combined approaches . . . . .	15
3.3	Summary . . . . .	15
<b>4</b>	<b>Conceptual Design and Implementation</b>	<b>16</b>
4.1	Preprocessing . . . . .	16
4.1.1	Vertical Byte Reading . . . . .	16
4.1.1.1	First Ideas . . . . .	16
4.1.1.2	Other Difficulties . . . . .	16

---

4.1.1.3	Compression improvements through vertical reading . . . . .	17
4.1.2	Varying maximum run lengths . . . . .	17
4.1.3	Byte remapping . . . . .	18
4.1.4	Burrows-Wheeler-Transformation Appliance . . . . .	19
4.2	Other Processing Options . . . . .	21
4.2.1	Performance improvements . . . . .	21
4.3	Implementation Decisions . . . . .	21
4.4	Implementation Detail . . . . .	21
4.4.1	Burrows Wheeler Transformation . . . . .	21
4.4.2	Byte Remapping . . . . .	21
4.5	Implementation Evaluation . . . . .	21
4.6	Summary . . . . .	21
<b>5</b>	<b>Evaluation</b>	<b>22</b>
5.1	Functional Evaluation . . . . .	22
5.2	Benchmarks . . . . .	22
5.2.1	Burrows-Wheeler-Transformation . . . . .	22
5.2.2	Vertical encoding . . . . .	22
5.3	Conclusion . . . . .	22
<b>6</b>	<b>Discussion</b>	<b>23</b>
	<b>Bibliography</b>	<b>24</b>

# List of Figures

2.1	Simple Example Huffman Tree with 3 Leaf Nodes . . . . .	7
4.1	Different run lengths for byte mapping and varying maximum run lengths . . . . .	19

# List of Tables

2.1	Entropy in relation to the Probability of Symbols . . . . .	4
2.2	T4 static Huffman codes . . . . .	6
2.3	Burrows Wheeler Transformation Matrix (all cyclic rotations) . . . .	8
2.4	State of the Art compression ratios . . . . .	9
3.1	Binary RLE on the Calgary Corpus . . . . .	10
3.2	Byte-wise RLE on the Calgary Corpus . . . . .	11
3.3	The file <i>pic</i> with increasing bits per RLE encoded number . . . . .	11
4.1	Binary RLE on vertical interpreted data . . . . .	17
4.2	Calgary Corpus encoded with vertical reading and 2 bits per RLE run	17
4.3	Calgary Corpus encoded with vertical reading and 2 bits per RLE run and 5 for higher significance bits . . . . .	18
4.4	Calgary Corpus encoded with vertical reading, byte remapping and 2 bits per RLE run and 5 for the 3 most significance bits . . . . .	19
4.5	Simple Burrows Wheeler Transformation (with 0x02 & 0x03) on byte wise RLE . . . . .	20
4.6	Burrows Wheeler Transformation on byte wise RLE . . . . .	20

# 1. Introduction

Die Einleitung besteht aus der Motivation, der Problemstellung, der Zielsetzung und einem ersten Überblick über den Aufbau der Arbeit.

TODO:

- explain compression ratio
- define unit of compression to quantify question and results

## 1.1 Motivation

In the last decades, digital data transfer became available everywhere and to everyone. This rise of digital data urges the need for data compression techniques or improvements on existing ones. Run-length encoding [10] (abbreviated as RLE) is a simple coding scheme that performs lossless data compression. RLE compression simply represents the consecutive, identical symbols of a string with a run, usually denoted by  $\sigma^i$ , where  $\sigma$  is an alphabet symbol and  $i$  is its number of repetitions. To give an example, the string `aaaabbaaabbba` can be compressed into RLE format as  $a^4b^2a^3b^4a^1$ . Its simplicity and efficiency make run-length encoding still usable in several areas like fax transmission, where RLE compression is combined with other techniques into Modified Huffman Coding [5]. Most fax documents are typically simple texts on a white background, RLE compression is particularly suitable for fax and often achieves good compression ratios. Another appliance of RLE is optical character recognition, in which the inputs are usually images of large scales of identically valued pixels [1].

## 1.2 Problem statement

Some strings like `aaaabbbb` achieve a very good compression rate because the string only has two different characters and they repeat more than twice. Therefore it can be compressed to  $a^4b^4$  so from 8 bytes down to 4 bytes if you encode it properly. On the other hand, if the input is highly mixed characters with few or no repetitions at all like `abababab`, the run length encoding of the string is  $a^1b^1c^1d^1e^1f^1g^1h^1i^1j^1$  which needs up to 16 bytes depending on the implementation.

So the inherent problem with run length encoding is obviously the possible explosion in size, due to missing repetitions in the input string. Expanding the string to twice the original size is not really a good compression so one has to make sure the input data is fitted for RLE as compression scheme. One goal is to improve the compression ratio on data suited for run length encoding and perform better than the originally proposed RLE. Another goal should be to minimize the increase in size in the worst case scenario.

## 1.3 Main Objective

Was ist das Ziel der Arbeit. Wie soll das Problem gelöst werden?

- bessere kompressionsrate im vergleich zu konventionellem RLE
- gleiche oder ähnliche decoding zeit ?
- ansatz beschreiben

The main objectives that derives from the problem statement is to archive an improved compression ratio compared to regular run length encoding. To unify the measurements, the compression ratio is calculated by encoding all files listed in the Galgary corpus and then normalize the results.

Since most improvements like permutations on the input, for example a revertable Burros-Wheeler-Transformation to increase the number of consecutive symbols or a different way of reading the byte stream take quite some time, encoding speed will increase. A second objective might be to keep decoding speed close to the original run length encoding.

## 1.4 Structure of this work

Was enthalten die weiteren Kapitel? Wie ist die Arbeit aufgebaut? Welche Methodik wird verfolgt?

- describe following structure
- use references
- try to keep idea of a recurrent theme



## 2. Principles of compression

The basic idea of compression is to remove redundancy in data. Compression can be broken down into two broad categories: Lossless and lossy compression. Lossless compression makes it possible to reproduce the original data exactly while lossy compression allows the some degradation in the encoded data to gain even higher compression at the cost of some of the original information. To understand compression, one first has to understand some basic principles of information theory like Entropy and different approaches to compress different types of data with different encoding. We will also show the key differences between probability coding and dictionary coding.

### 2.1 Compression and Encoding fundamentals

TBD

...

#### 2.1.1 Information Theory and Entropy

As Shannon described his analysis about the English language [11], he used the term Entropy closely aligned with its Definition in classical physics, where it is defined as the disorder of a system. Specifically speaking, it is assumed that a system has a set of states it can be in and it exists a probability distribution over those states. Shannon then defines the Entropy as:

$$H(S) = \sum_{s \in S} P(s) i(s)$$

where  $S$  describes all possible States,  $P(s)$  is the likelihood of  $s \in S$ . So generally speaking it means that evenly distributed probabilities imply a higher Entropy and vice versa. It also implies that, given a source of information  $S$ , its average information content per message from  $S$  is also described by this formula.

In addition to that, Shannon defined the term self information  $i(s)$  as:

$$i(s) = \log_2 \frac{1}{P(s)}$$

indicating that the higher the probability of a state, less information can be contained. As an example, the statement “The criminal is smaller than 2 meters.” is very likely but doesn’t contain much information, whereas the statement “The criminal is larger than 2 meters.” is not very likely but contains lots of information. With these definitions in mind, we can analyze some properties for the English language from an information theory perspective of view.

	$P(a)$	$P(b)$	$P(c)$	$P(d)$	$P(e)$	$H$
1.	0.2	0.2	0.2	0.2	0.2	2.322
2.	0.94	0.01	0.01	0.01	0.01	0.322

Table 2.1: Entropy in relation to the Probability of Symbols

Combining those approaches, you will find that  $\Sigma$  being a finite alphabet, and  $P(s_i)$  describing the likelihood of  $s_i$  and  $i(s_i)$  its self information,  $i(s_i)$  also describes the amount of bits needed for this symbol, therefore its  $\log_2$ .

$$H(\Sigma) = \sum_{i=1}^n P(s_i) \cdot \log_2 \frac{1}{P(s_i)}$$

So the entropy also describes the theoretical amount of bits needed to persist a message from  $\Sigma$ , calculated in *bps* as Bits per Symbol.

TODO: describe some relations between probability and entropy of information source

### 2.1.2 General Analysis

To evaluate the efficiency of a specific compression technique, we have to determine how much information the raw data contains. In this case for textual compression at first, we are talking about the English language. There have been broad analysis of ASCII Entropy, consisting of 96 different printable symbols for the English language, generating approaches for calculating the entropy [8].

If we assume a 96 Symbol alphabet and a uniform probability distribution, we have a quite high entropy of  $\log_2(96) = 6,6bps$ . In a empirically distribution generated by text analysis, the entropy comes down to  $4.5bps$ . Using an encoding with separated encoding per Symbol like the Huffman Coding, we can archive an entropy of  $4.7bps$  which is only slightly worse than the theoretical entropy. By changing the assumption to blocks of symbols of length 8, we get  $96^8$  different blocks. With a probability distribution close to the English language we get an entropy as low as  $1.3bps$  which implies a possible reduction of symbols to 3 printable characters without generating longer texts.

Consulting newer sources we will find that, up to 500 character long symbol based text analysis, with around 20.3 Mio. different symbols, results in an entropy of around  $1,58bps$  [4]. This gives a clear limit of how much compression we can theoretically expect from English text. Using a combination of different approaches like Huffman Coding and Run Length Coding, and encoding different bits of a single symbol with different approaches, we do no longer encode every symbol separately, which might result in a better compression.

### 2.1.3 Probabilistic Coding

The general idea behind Probability Coding is to analyze probabilities for messages and the encode them in bit strings according to their probability. This way, messages that are more likely and will repeat more often can be encoded in a smaller bit string representation. The generation of those probability distributions is considered part

of the Analyzer module by the algorithm and will be discussed later on in Section 4. Probability coding compression can be discerned into fixed unique coding, which will represent each message with a bit string with the amount of bits being an integer, for example Huffman coding as type of prefix coding. In contrast to Huffman coding there are also arithmetic codes which can represent a message with a floating point number of bits in the corresponding bit string. By doing so they can “fuse” messages together and need less space to represent a set of messages.

### 2.1.4 Dictionary Coding

Dictionary Coding is best suited for Data with a small amount of repeating patterns like a text representing specific words. In this case it is very effective to save the patterns just once and refer to them if they are used later on in the text. If we know a lot about the text itself in advance, a static dictionary is sufficient but in general we want to compress an arbitrary text. This requires a more general approach as used by the well known LZ77 and LZ88 algorithms described by Jacob Ziv and Abraham Lempel [12]. They use a so called sliding window principle, a buffer moving across the data to encode in the case of LZ77 or a dynamic dictionary in the implementation of LZ78. Both of them are still used in modified versions in real world applications as described in Section 2.5.

The main difference between probabilistic coding and dictionary coding is, that the so far presented probability based methods like RLE or Huffman Coding are working on single character or byte whereas the dictionary based methods encode groups of characters of varying length. This is a clear advantage in terms of compression performance on every day data because of its repeating patterns most textual data has, as shown in Section 3.1.

### 2.1.5 Irreversible Compression

Irreversible Compression or also called “Lossy Compression” is a type of compression which loses information in the compression process and is therefore not completely reversible, hence the name. There are a lot of use cases for these type of compression algorithms, mostly used for images, video and audio files. These files typically contain information almost not perceptible for an average human like really small color differences between two pixels of an image or very high frequencies in an audio file. By using approximations to store the information and accept the loss of some information while retaining most of it, lossy image compression algorithms can get twice the compression performance compared to lossless algorithms. Due to the fact that most lossy compression techniques are not suited for text compression which is the main use case for this work, we will not elaborate any further on this topic.

## 2.2 Run Length Coding

Run length coding might count as the simplest coding scheme that makes use of context and repeating occurrences as described in the Section 1.1. While the example made earlier was in textual representation, it is best suited and mostly used for pallet based bitmap images [3] such as fax transmissions or computer icons.

run length	white codeword	black codeword
0	00110101	0000110111
1	000111	010
2	0111	11
3	1000	10
4	1011	011
...		
20	0001000	00001101000
...		
64+	11011	0000001111
128+	10010	000011001000

Table 2.2: T4 static Huffman codes

### 2.2.1 The history

The ITU-T T4 (Group 3) standard for Facsimile (fax) machines [6] is still in force for all devices used over regular phone lines. Each transmission sends a black and white image, where each pixel is called a *pel* and with a horizontal resolution of  $8.05 \frac{pels}{mm}$  and the vertical resolution depending on the mode. To encode each sequence of black and white pixels, the T4 standard uses RLE to encode each sequence of black and white pixels and since there are only two values, only the run length itself has to be encoded. It is assumed that the first run is always a white run, so there is a dummy white pel at the beginning of each sequence. For example, the sequence *bbbbwubbbb* can be encoded as 1,4,2,5 with the leading white dummy pixel.

To further reduce the RLE encoded sequence, a probability coding procedure like Huffman coding can be used to code the sequence, since shorter runs like lengths are generally more common than very long runs, which is also done by the T4 standard. Based on a broad analysis of these run lengths, the T4 defined a set of static Huffman codes, using a different code for the black and white pixels.

A small subset of the static T4 Huffman codes are given in Table 2.2.1 which also shows the most likely runs to occur. Runs of more than 64 have to be encoded in multiple Huffman codes, for example a run of 150 has to be encoded with the Huffman code of 128 followed by the code for 22.

By combining RLE with a probabilistic approach, the ratio of compression rises because we no longer have to encode a run of length 1 or 2 with a fixed size of bits each instead we can write a recognizable Huffman code of varying length, which will be explored in greater detail later on as well as the decoding of the shown Huffman codes.

### 2.2.2 Limitations

As mentioned in the Section 1.2, run length encoding is rarely used for regular text or continuous tone images because its potentially increase in size, due to non repetitive characters or bytes. To reduce this problem there are several approaches known, partly implemented and used some of which will be addressed like a Burrows-Wheeler-Transformation.

### 2.2.3 Run Length Encoding today

While it is still in use by fax transmission or other highly specific tasks, it is mostly used in combination with other approaches. TODO: give more references [7]

## 2.3 Prefix Coding

Prefix codes make use of an untypical idea in computer science. Usually we deal with fixed length codes like 7 bit ASCII or 32 bit Integer representations which map each possible value into that fixed amount of bits. For compression it would be a benefit if we could write codes of variable length. The problem with these varying length codes is that as soon as they are part of a sequence, it becomes very hard to tell where one code starts and finishes, resulting in ambiguous encoding. For example the set of codes  $\{(a, 1), (b, 01), (c, 101), (d, 011)\}$  and the encoded message is 1011, there is no way to tell which message was encoded.

One way to address this problem is by adding extra stop symbols or encoding a length before each code but those just add additional data. Another approach is to use so called prefix codes and make no code which is prefix of another code, which makes them uniquely decodable as show in the next section.

### 2.3.1 Huffman Coding

Huffman coding

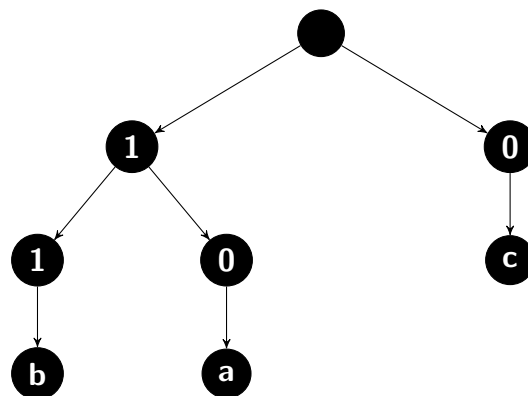


Figure 2.1: Simple Example Huffman Tree with 3 Leaf Nodes

TODO: explain stuff ...

## 2.4 Other methods

There are other well known methods like

- move to front
- arithmetic coding
- real world applications are mostly combinations of those

...

### 2.4.1 Burrows-Wheeler-Transformation

The Burrows-Wheeler-Transformation is not a compression technique but rather a method to prepare text to increase its compression potential, described by M. Burrows and D. J. Wheeler [2]. It is a Transformation of a string  $S$  of  $n$  characters by forming the  $n$  rotations of  $S$ , sorting them lexicographical, and extracting the last character of each of the rotations. The result of this transformation is a string  $L$ , consisting of these last characters of each sorted rotation. The algorithm also has to add special termination symbols or compute the index  $I$  of the original string  $S$  in the sorted list of rotations. Surprisingly, there is an efficient algorithm to invert the transformation back to the original string  $S$  given only  $L$  and the Index  $I$  [9].

As an example of the creation of  $L$ , given the input string  $S = \text{'abraca'}$ ,  $n = 6$ , and the alphabet  $X = \{a, b, c, r\}$ . Create a  $N \times N$  matrix  $M$  whose elements are characters, and whose rows are rotations of  $S$ . Sort all rotations in lexicographical order.

In this example, the index is  $I = 1$  and the matrix  $M$  is

row 0	a	a	b	r	a	c
row 1	a	b	r	a	c	a
row 2	a	c	a	a	b	r
row 3	b	r	a	c	a	a
row 4	c	a	a	b	r	a
row 5	r	a	c	a	a	b

Table 2.3: Burrows Wheeler Transformation Matrix (all cyclic rotations)

The resulting string  $L$  corresponds to the last column of  $M$ , with characters  $M[0, n-1]$ ,  $\dots$ ,  $M[n-1, n-1]$ . The output of the transformation is the pair  $(L, I)$ , in the example,  $L = \text{'caraab'}$  and  $I = 1$ . Obviously the string  $L$  contains consecutive identical characters, which result in better compressibility as shown in 3.2.1.

### 2.4.2 Inverse Burrows-Wheeler-Transformation

TODO: explain reversing a bwt

## 2.5 State of the Art

State of the art compression

- techniques
- use cases
- limits

## 2.6 Limits of compression

- existence of not compressable strings
- <https://www.quora.com/Is-there-a-theoretical-limit-to-data-compression-If-so-how-was-it-found>
- unable to compress random data
- Kolmogorow-Komplexität

method	size in bytes	compression	ratio in $\frac{bits}{symbol}$
uncompressed	3,141,622	100%	8.00
compress	1,272,772	40.0%	3.24
ZIP v2.32	1,020,781	32.4%	2.59
gzip v1.3.5	1,017,624	32.3%	2.58
bzip2 v9.12b	828,347	26.3%	2.10
ppmd	740,737	23.5%	1.88
ZPAQ v7.15	659,709	20.9%	1.67

Table 2.4: State of the Art compression ratios

## 3. Analysis

The following chapter contains a detailed analysis of the problem and some fundamental requirements for the algorithm. To have a comparison to some extent, the initial performance analysis was performed on the Calgary Corpus but the results of the unmodified run length coding algorithm were very underwhelming as expected.

### 3.1 Initial Findings

Originally developed and used on black and white pallet images containing only two values, on a binary level so to speak, we want to use it for rather arbitrary data, mostly text. The initial algorithm is not suited for that because continuous text as binary representation does not contain runs of any kind, which could be compressed. The ASCII representation of the letter 'e' which is the most common in general English, has the value 130 or '01100101' as 7-bit ASCII or '001100101' as byte value of the UTF-8 representation. As you can see, there are no runs of a considerable size and this is the case for most printable characters as they all have a value between 32 and 127 (or 255 for the extended ASCII). Applied to the Calgary Corpus in this simple implementation, there should be an increase in size expected or *negative compression* as one might say.

bits per rle number	expansion ratio	bits per symbol in $\frac{bits}{symbol}$
8	3.298046423741734	26.38437138993387
7	2.8897459975751163	23.11796798060093
6	2.4839251325134675	19.87140106010774
5	2.0825541259578895	16.66043300766311
4	1.6895640359371056	13.516512287496845
3	1.3130541262757818	10.504433010206254
2	1.04555716691706	8.36445733533648

Table 3.1: Binary RLE on the Calgary Corpus

The results in Table 3.1 depict the anticipated, an increase in size regardless of the amount of bits used to encode a run. By using 8 bits to encode a single run, in the worst case scenario a byte which is only alternating values like 01010101, would expand to 8 bytes, all encoding a run of length 1. For this reason, many implementations combine RLE with other encoding schemes like Huffman encoding to be able to encode runs with variable length.

RLE is also applicable on a byte level, because there should be repetitions of any kind like consecutive letters or line endings (EOL). Byte level RLE encodes runs



of identical byte values, ignoring individual bits and word boundaries. The most common byte level RLE scheme encodes runs of bytes into 2-byte packets. The first byte contains the run count of 1 to 256, and the second byte contains the value of the byte run. If a run exceeds a count of 256, it has to be encoded twice, one with count 256 and one with any further runs.

bits per rle number	ratio in %	bits per symbol in $\frac{bits}{symbol}$
8	165	13.2
7	154	12.38
6	144	11.57
5	134	10.77
4	125	10.00
3	116	9.29
2	109	8.74

Table 3.2: Byte-wise RLE on the Calgary Corpus

However after some analysis of the corpus data, it was shown that most runs had a value of one and almost no runs larger than 4 occurred, which lead to the conclusion, two bit for the run count should be plenty, which is also shown in Table 3.1. Even with a run size of just two bits, there is still a increase in size of about 9% and uses  $8.74 \frac{bits}{symbol}$ . This is still useful as a kind of a base line. Interestingly the binary implementation performs better on 2 bits per RLE number (4 % increase in size) than the byte implementation (9 % increase in size) but also worse with a higher amount of bits per run, where it expands the data to more than triple in size. It is unclear which kind of implementation will profit most of preprocessing, so both will be further analyzed, but the benefit of byte wise RLE is the better worst case performance of 1.5 up to 2 times the original size compared to binary RLE with up to 4.5 times the original size.

If we take a more detailed look, we can see that while most files expand with larger RLE numbers, some files have their minimum size when encoded with higher RLE numbers of up to 7 bit. With the simple binary based RLE, almost all files of the Calgary Corpus expand linear related to the amount of bits used for the encoding. However the file *pic* decreases in size until 7 bits per RLE number used to a sizes of just 19.5% of its original size with only  $1.56 \frac{bits}{symbol}$  while the other files just doubled or even tripled in size. Using the byte wise operating RLE we see a similar result with the file *pic*, but not as decent with 27.2% of its original size using  $2.17 \frac{bits}{symbol}$  encoding with 6 bits per run.

file	size original	$\frac{bits}{RLEnumber}$	size encoded	ratio in %	$\frac{bits}{symbol}$
pic	513216	2	350292	68.25	5.46
		3	235067	45.80	3.66
		4	165745	32.29	2.58
		5	126349	24.61	1.96
		6	106773	20.80	1.66
		7	100098	19.50	1.56
		8	101014	19.68	1.57

Table 3.3: The file *pic* with increasing bits per RLE encoded number

An additional step of the improvement could be the detection of high efficiency with regular binary RLE to keep files highly suited for this method at its peak.

## 3.2 Possible Improvements by Preprocessing

The broad idea of preprocessing is to manipulate the input data in a way that results in data which can be compressed more efficiently than the original data. This can be done in various ways, some of them will be explored in greater detail to find out if it is implementable or not. One way of doing so is a Burrows-Wheeler-Transformation.

### 3.2.1 Burrows-Wheeler-Transformation

To understand how a Burrows-Wheeler-Transformation improves the effectiveness of compression, consider the effect in a common word in English text. Examine the letter ‘t’ in the word ‘the’, in an input string holding multiple instances of this word. Sorting all rotations of a string results in all rotations starting with ‘he ’ will be sorted together and most of them are going to end in the letter ‘t’. This implies that the transformed string L has a large number of the letter t, combined with some other characters, such as space, ‘s’, ‘T’, and ‘S’. This is true for all characters, so any substring of L is likely to contain a large number of a some distinct characters. “The overall effect is that the probability that given character  $c$  will occur at a given point in L is very high if  $c$  occurs near that point in L, and is low otherwise.” [2]

It is obvious that this should always improve the performance of byte level RLE because the transformation is taking place at character level but it should not effect the binary implementations.

### 3.2.2 Vertical byte reading

Instead of performing compute intense operation on the data, we could also interpret the data in a different way and apply the original run length encoding on binary data. This idea is also know for binary RLE on images, where the encoding in the image follows a specific path.

TODO showcase

By reading the data in chunks of a fixed size, it is possible to read all most significant bits of all bytes, then the second most significant bits of all bytes and so on. This interpretation results in longer runs as shown in the example below.

The binary UTF-8 interpretation of the example string from earlier  $S = \text{‘abraca’}$  results in 8 runs of length 1, 9 runs of length 2 as well as 3 runs of length 3 and 4.

48 ELEMENTS: \_\_\_\_\_  
 01110000010011000010001110001000110000010110000110  
 1110000011

Reading the data in a different way, all most significant bits, then al second most significant bits and so forth, it results in much longer runs as shown below.

2-D: \_\_\_\_\_  
 01110000011

<Row 1>

---

01100010	<Row 2>
011110010	<Row 3>
01100001	<Row 4>
01100011	<Row 5>
01100001	<Row 6>

---

Now we have 5 runs of length 6, 2 runs of length 3, 3 runs of length 2 and just 6 runs of length 1 as opposed by the simple interpretation. This is because the binary similarity between the used characters, as the character for a and b only differ in one bit. It is clear that simply a different way of reading the input does not compress the actual data, instead it enables a better application of existing compression.

### 3.2.3 Byte remapping

The effect of very long runs in the last example was mainly because the binary representations of the used characters are very similar, so the range of byte values used was very small (between a = 97 and r = 114). Introducing other used symbols like uppercase letters, space or new lines, the used range expands.

The binary representation of a String like S' = "Lorem ipsum dolor sit amet, consectetur adipiscing elit.", results in a worse result as shown below. The usage of other characters expanded the used byte range to between 32 and 117 which results in shorter average runs. Interestingly the most significant bit is always 0, a fragment from the backwards compatibility with standard ASCII encoding.

2-D: 

---

01001100	<Row 1>
01100111	<Row 2>
011110010	<Row 3>
0110001	<Row 4>
01101101	<Row 5>
00010000	<Row 6>
01101001	<Row 7>
01111000	<Row 8>
01111001	<Row 9>
0111101	<Row 10>
01101101	<Row 11>
00010000	<Row 12>
01100010	<Row 13>
01101111	<Row 14>
01101100	<Row 15>
01101111	<Row 16>
011110010	<Row 17>
00010000	<Row 18>
01111001	<Row 19>
01101001	<Row 20>
01111010	<Row 21>
00010000	<Row 22>
01100001	<Row 23>
01101101	<Row 24>

---

0111001011	<Row 25>
0111110100	<Row 26>
0011011110	<Row 27>

---

One idea to solve the shorter runs might be a dynamic byte remapping, as the input data is read in parts, where the most frequently used bytes are mapped to the lowest value. This way the values are not alternating in the whole range of 0 to 255 but rather in a smaller subset and the most frequent ones will be the smallest values, so in theory our average runs should increase because we should encounter more consecutive zeros. Some sections have more specific characters or bytes than others but this idea can also be applied to the whole file however it is unclear at this point what method outperforms which. A single map for each block of data should result in lower average values used but also creates a kind of overhead because the mapping has to be stored in the encoded file as well. Applying a simple mapping to lower values results in the following horizontally interpreted rows.

2-D:

---

00001001	<Row 1>
000000010	<Row 2>
000000101	<Row 3>
000000011	<Row 4>
000000001	<Row 5>
000000000	<Row 6>
000000100	<Row 7>
000001101	<Row 8>
000000110	<Row 9>
000001110	<Row 10>
000000001	<Row 11>
000000000	<Row 12>
000001011	<Row 13>
000000010	<Row 14>
000001100	<Row 15>
000000010	<Row 16>
000000101	<Row 17>
000000000	<Row 18>
0000000110	<Row 19>
0000000100	<Row 20>
0000000111	<Row 21>
0000000000	<Row 22>
000001010	<Row 23>
0000000001	<Row 24>
0000000011	<Row 25>
0000000111	<Row 26>
000001000	<Row 27>

---

Using this method, the 4 most significant bits all result in zero columns, even row 5 has long runs while it is worth noting that the mapping itself has to be persisted in

the encoded file as well. It is also still unclear if this idea scales well or is applicable to other files.

### 3.2.4 Combined approaches

The idea of combining different compression methods into a superior method is not new and was also performed on RLE as mentioned in Section 2.2.1. While the idea of encoding the RLE numbers with Huffman codes is already known and analyzed, it is mostly in a static sense and optimized for special purpose applications. However the vertical byte reading enables new approaches, even more in combination with the idea of byte remapping and might become applicable to more than just binary Fax transmissions or DNA sequencing with longer runs of any kind in average.

It might be interesting to see how well the appliance performs using the vertical binary encoding, in combination with the byte mapping. We expect the more significant the bit, the longer the runs because alternating values should be mostly on the lower significance bits. Using larger run length numbers for these rows while using smaller RLE numbers or even another encoding scheme like simple Huffman encoding we should improve our initial results.

Another interesting approach is the improvement archived by the combination of different methods like performing this modified RLE run on arbitrary data and then encode the results with Huffman codes, using shorter codes for more frequent runs. This idea is not new and also used in the current Fax Transmission Protocol but only for the simple binary RLE in combination with modified Huffman Codes. Other papers also mentioned these combined approaches and seemed to archived good compression ratios, not much worse than the theoretical limit of around  $1.5 \frac{\text{bits}}{\text{symbol}}$  shown in Section 2.1. This was also done by M. Burrows and D.J. Wheeler in 1994 with their Transformation, in combination with a Move-to-Front Coder and a Huffman Coder [2]. Encoding the Calgary Corpus resulted in a decrease in size to just 27% of its original with a mean  $\frac{\text{bits}}{\text{symbol}}$  of just 2.43. This approach would no longer be considered preprocessing, but if more compression could be archived by adding a post-processing step it is still worth trying out.

## 3.3 Summary

TODO

Am Ende sollten ggf. die wichtigsten Ergebnisse nochmal in *einem* kurzen Absatz zusammengefasst werden.

## 4. Conceptual Design and Implementation

Without further ado, design and implementation of the initial ideas began. As Section 3 showed, there are some potentially promising improvements to be made to run length encoding, but how well they scale and work on a larger input with versatile symbols or bytes has to be determined.

### 4.1 Preprocessing

#### 4.1.1 Vertical Byte Reading

Implementing the vertical reading of the input shown in Section 4.1.1 was not hard but it should be kept in mind that the size of the input chunks has to be divisible by the size of 8. Otherwise parsing it into an Array of Bytes results in the last Byte having some padding which might cause problems later on. By collecting the bits into a proprietary data structure, we avoid this problem but working with bytes internally should be easier. For larger files it should also be possible to work in space and writing the output on the fly during the reading process. Nonetheless both ways of collecting all bits of identical significance and writing on the fly are implemented.

##### 4.1.1.1 First Ideas

Initially some other ideas have been followed with very poor results. One idea was arranging all input bits in a Matrix or square Matrix in a way it would still be receivable later on. This way other methods from linear algebra would have been applicable to the input data, so that the construction of a triangular matrix or other preprocessing would result in long runs of zeros. Difficulties in the construction and transformation of the Data lead to the abandonment of this approach and the already described vertical interpretation was used further on.

##### 4.1.1.2 Other Difficulties

Other difficulties arose while performing bit operations, because it became slow on larger workloads.

TODO: explain used lib

#### 4.1.1.3 Compression improvements trough vertical reading

First results of just plain binary RLE on the vertical interpretation improved its overall performance and archived a small edge over regular RLE with a reduction in the corpus size of around 3% with 2 bits per run.

bits per rle number	ratio in %	bits per symbol in $\frac{bits}{symbol}$
3	110.82585915202825	8.86606873216226
2	97.14322771462668	7.771458217170134

Table 4.1: Binary RLE on vertical interpreted data

If we take a closer look on each file, we see similar results compared to the original proposed binary RLE, where most files had a compression ratio of around 1 with 2 bits per RLE encoded number, except for the file *pic*. Average sizes are increasing again with more bits per run up to 2.5 times its original size and also the file *pic* has the best compression ratio. This time although it is at its peak using 5 bits per RLE run and only archives a compression of  $3.67 \frac{bits}{symbol}$  compared to  $1.56 \frac{bits}{symbol}$  with simple binary RLE.

file	size original	size encoded	ratio in %	$\frac{bits}{symbol}$
bib	111261	117287	105.41	8.43
book1	768771	781773	101.69	8.13
book2	610856	621173	101.68	8.13
geo	102400	105454	102.98	8.23
news	377109	377779	100.17	8.01
obj1	21504	22003	102.32	8.18
obj2	246814	265075	107.39	8.59
paper1	53161	54118	101.80	8.14
paper2	82199	83467	101.54	8.12
pic	513216	375791	73.22	5.85
progc	39611	38896	98.19	7.85
progl	71646	68428	95.50	7.64
progp	49379	48629	98.48	7.87
trans	93695	91883	98.06	7.84
all files	3145718	3055852	97.14	7.77

Table 4.2: Calgary Corpus encoded with vertical reading and 2 bits per RLE run

These results were already better than the simple binary RLE but they were still quite far from the expected outcome. This mainly arose because increasing the bits per run improved the result for the most significant bits but also degraded the results for the other bits. So the idea of encoding the bits of different significance with different RLE schemes with varying bits per run.

#### 4.1.2 Varying maximum run lengths

This time the most significant bits have been encoded with more bits per run than the other ones and after bench marking every combination, it turned out, with 2

bits per run and 5 bits per run for higher bits, it improved by another 4 percent. Now most files have a  $\frac{\text{bits}}{\text{symbol}}$  ratio below 8, more specifically all textual files, which relates to the earlier mentioned ASCII fragments in UTF-8 encoding. This time higher possible runs on this position resulted in fewer runs in total, so in a better compression overall.

file	size original	size encoded	ratio in %	$\frac{\text{bits}}{\text{symbol}}$
bib	111261	110258	99.09	7.92
book1	768771	733208	95.37	7.62
book2	610856	582583	95.37	7.62
geo	102400	121496	118.64	9.49
news	377109	353955	93.86	7.50
obj1	21504	24067	111.91	8.95
obj2	246814	292143	118.36	9.46
paper1	53161	50758	95.47	7.63
paper2	82199	78273	95.22	7.61
pic	513216	356712	69.50	5.56
progc	39611	36392	91.87	7.34
progl	71646	63902	89.19	7.13
progp	49379	45508	92.16	7.37
trans	93695	85964	91.74	7.33
all files	3145718	2939315	93.43	7.47

Table 4.3: Calgary Corpus encoded with vertical reading and 2 bits per RLE run and 5 for higher significance bits

Applying this increase to more than the most significant bit lead to a decrease in performance, which most likely related to the shorter runs on lower order bits, seen in Section 3.2.2. Therefore the concept of mapping the input to lower value bytes to artificially increase runs of consecutive zeros on more significant bits had to be applied first.

### 4.1.3 Byte remapping

As shown in Section 3.2.3 this effect could become useful if it resulted in longer runs. This effect was also seen in the second and third most significant bit, as higher value bytes became unlikely in the input data after the remapping so the idea of varying maximum run lengths for different significant bits became more appealing again. To determine the best combination of maximum run lengths and how many most significant bits should be encoded with a higher maximum run length, all combinations of these were tested and the results plotted below.

TODO: fix plot scaling



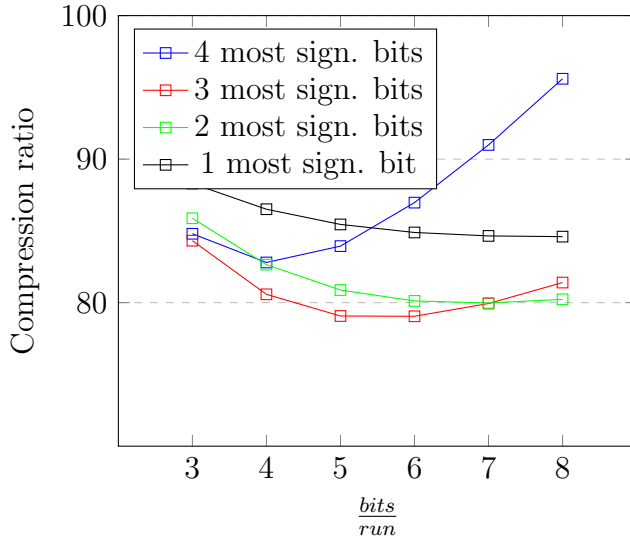


Figure 4.1: Different run lengths for byte mapping and varying maximum run lengths

Empirically, the combination of 2 bits per run and 5 bits per run for the 3 most significant bits yielded the overall best results with 79% of its original size and  $6.32 \frac{\text{bits}}{\text{symbol}}$ . It has to be mentioned, that the expected overhead is rather small but the increase in average run length lead to a reduction in size.

file	size original	size encoded	ratio in %	$\frac{\text{bits}}{\text{symbol}}$
bib	111261	97867	87.96	7.03
book1	768771	579649	75.39	6.03
book2	610856	484989	79.39	6.35
geo	102400	116487	113.75	9.10
news	377109	315446	83.64	6.69
obj1	21504	23871	111.00	8.88
obj2	246814	276301	111.94	8.95
paper1	53161	43556	81.93	6.55
paper2	82199	62544	76.08	6.08
pic	513216	272034	53.00	4.24
prog	39611	33073	83.49	6.67
progl	71646	53653	74.88	5.99
progp	49379	39400	79.79	6.38
trans	93695	84494	90.17	7.21
all files	3145718	2487460	79.07	6.32

Table 4.4: Calgary Corpus encoded with vertical reading, byte remapping and 2 bits per RLE run and 5 for the 3 most significance bits

#### 4.1.4 Burrows-Wheeler-Transformation Appliance

Another possible preprocessing step which promised an improvement is the mentioned Burrows-Wheeler-Transformation from Section 2.4.1, initially applied to regular binary and byte wise RLE. By mistake a very simple transformation implementation was chosen, working by adding additional start and stop symbols to the

input string (0x02 as STX, start of text and 0x03 as ETX, end of text). Some basic testing and playing around worked great but later on it revealed some major issues. For example the Calgary Corpus consists of more than textual data, in fact the files geo, obj1, obj2 and pic contain of some binary data of include the symbols STX or ETX so we wont be able to apply the transformation to these. Another shortcoming was the very poor time complexity of almost  $O(n^2)$  because under the hood, it uses a dual pivot Quick-sort algorithm from the JDK 11, which is typically faster than traditional one pivot Quick-sort. This algorithm offers  $\Theta(n \log(n))$  average time complexity but in the worst case, its time complexity is cubic. This problem was partially solved by reading the input data in parts and performing the transformation on each part, result in a much smaller length  $n$  and thus better run time at the expense of a slightly worse transformation result. As all chunks are individual transformations, they can also be computed in parallel without much effort.

bits per rle number	ratio in %	bits per symbol in $\frac{bits}{symbol}$
3	95.4169856525027	7.633358852200216
2	91.391	7.311309412577118

Table 4.5: Simple Burrows Wheeler Transformation (with 0x02 & 0x03) on byte wise RLE

While it was only applicable to textual data and very slow, even when divided into smaller parts and computed in parallel, it improved the overall results of byte wise RLE by 16% to a compression ratio of slightly over 7  $\frac{bits}{symbol}$  which seemed like a good start. Regular binary RLE did not really benefit from this transformation as expected but on vertical interpretation, consecutive characters result in successive bits on every significance. Still this implementation had to be dropped and switched against one that could handy arbitrary input to be able to transform all files. This time all files could be processed and the resulting compression with byte wise RLE improved further.

bits per rle number	ratio in %	bits per symbol in $\frac{bits}{symbol}$
3	91.62970107301417	7.3303760858411335
2	89.46313687368035	7.157050949894428

Table 4.6: Burrows Wheeler Transformation on byte wise RLE

In Section 2.4 the Burrows-Wheeler-Transformation inversion was performed using a Matrix M containing all cyclic rotations of the input word, sorted in lexicographic order. However this is has a bad complexity and the inverting process does not even has access to this Matrix M.

TODO: explain creation using suffix array

TODO: create inversion using Lyndon words

## 4.2 Other Processing Options

- encoding with huffman codes
- ...

### 4.2.1 Performance improvements

## 4.3 Implementation Decisions

- why kotlin
- performance improvements with the graalvm
- other decisions
- ...

## 4.4 Implementation Detail

- detailed information about specific modules and classes
- ...

### 4.4.1 Burrows Wheeler Transformation

- TBD

### 4.4.2 Byte Remapping

- show use case

## 4.5 Implementation Evaluation

- evaluation of implementation choices made

## 4.6 Summary

Am Ende sollten ggf. die wichtigsten Ergebnisse nochmal in *einem* kurzen Absatz zusammengefasst werden.

# 5. Evaluation

Hier erfolgt der Nachweis, dass das in Kapitel 4 entworfene Konzept funktioniert. Leistungsmessungen einer Implementierung werden immer gerne gesehen.

## 5.1 Functional Evaluation

- Mathematical Comparison
- Comparison of encoded file sizes
- comparing to regular REL and Huffman coding
- ...

## 5.2 Benchmarks

- Benchmark with the Galgary corpus  
<http://www.data-compression.info/Corpora/>
- ...

### 5.2.1 Burrows-Wheeler-Transformation

bits per rle number	ratio in %	bits per symbol in $\frac{bits}{symbol}$
3	95.4169856525027	7.633358852200216
2	91.391	7.311309412577118

### 5.2.2 Vertical encoding

## 5.3 Conclusion

Am Ende sollten ggf. die wichtigsten Ergebnisse nochmal in *einem* kurzen Absatz zusammengefasst werden.

## 6. Discussion

(Keine Untergliederung mehr)

# Bibliography

- [1] Al-Okaily A, Almarri B, Al Yami S, and Huang CH. Toward a better compression for dna sequences using huffman encoding. *J Comput Biol.*, 24(4):280–288, 2017.
- [2] M. Burrows and D. J. Wheeler. A block-sorting lossless data compression algorithm. Technical report, 1994.
- [3] MJ Dallwitz. An introduction to computer images. *TDWG Newsletter*, 7(10):2, 1992.
- [4] Fabio G. Guerrero, editor. *A new look at the classical entropy of written English*, 2009.
- [5] Roy Hunter and A. Harry Robinson, editors. *International digital facsimile coding standards, Proceedings of the IEEE 68*, 1980.
- [6] International Telecommunication Union (ITU). *T.4 : Standardization of Group 3 facsimile terminals for document transmission*, 2004.
- [7] Péter Lehotay-Kéry and Attila Kiss. Genpress: A novel dictionary based method to compress dna data of various species. In Ngoc Thanh Nguyen, Ford Lumban Gaol, Tzung-Pei Hong, and Bogdan Trawiński, editors, *Intelligent Information and Database Systems*, pages 385–394, Cham, 2019. Springer International Publishing.
- [8] Maciej Liśkiewicz and Henning Fernau. *Datenkompression*.
- [9] Daisuke Okanohara and Kunihiro Sadakane. A linear-time burrows-wheeler transform using induced sorting. In Jussi Karlgren, Jorma Tarhio, and Heikki Hyvärinen, editors, *String Processing and Information Retrieval*, pages 90–101. Springer Berlin Heidelberg, 2009.
- [10] A. Harry Robinson and C. Cherry, editors. *Results of a prototype television bandwidth compression scheme, Proceedings of the IEEE 3*, volume 55, March 1967.
- [11] C. E. Shannon. *Prediction and entropy of printed English*, volume 30. Jan 1951.
- [12] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3):337–343, May 1977.

# Eidesstattliche Erklärung

Hiermit erkläre ich, dass ich diese Bachelor-/Masterarbeit selbständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt und die aus fremden Quellen direkt oder indirekt übernommenen Gedanken als solche kenntlich gemacht habe. Die Arbeit habe ich bisher keinem anderen Prüfungsamt in gleicher oder vergleichbarer Form vor-gelegt. Sie wurde bisher auch nicht veröffentlicht.

Trier, den xx. Monat 20xx