

# Seminar: Datenkompression

## Bzip2

Florian Brohm, 7443251      Toprak Saricerci, 7445073

Abgabe: 23. Februar 2023

# Inhaltsverzeichnis

<b>1</b>	<b>Was ist Bzip2?</b>	<b>4</b>
<b>2</b>	<b>Was ist eine Datenkompression?</b>	<b>4</b>
2.1	Verlustfreie und verlustbehaftete Kompressionen . . . . .	4
2.2	Entropiekodierungen und Substitutionskompressionen . . . . .	5
<b>3</b>	<b>Run-Length encoding</b>	<b>5</b>
3.1	Kodierung . . . . .	5
3.2	Dekodierung . . . . .	6
3.3	Effizienz . . . . .	6
3.4	Implementierung . . . . .	6
3.5	Pseudocode . . . . .	7
3.5.1	Kodierung . . . . .	7
3.5.2	Dekodierung . . . . .	7
<b>4</b>	<b>Huffman encoding</b>	<b>8</b>
4.1	Präfixfreiheit . . . . .	8
4.2	Kodierung . . . . .	8
4.3	Dekodierung . . . . .	9
4.4	Effizienz . . . . .	9
4.5	Pseudocode . . . . .	9
4.5.1	Tabellengenerierung . . . . .	9
4.5.2	Kodierung . . . . .	10
4.5.3	Dekodierung . . . . .	10
<b>5</b>	<b>Transformation</b>	<b>11</b>
5.1	Was ist eine Transformation? . . . . .	11
5.2	Warum sind Transformationen nützlich? . . . . .	11
<b>6</b>	<b>Move-to-front transform</b>	<b>12</b>
6.1	Erläuterung . . . . .	12
6.2	Inverse Transformation . . . . .	13
6.3	Implementierung . . . . .	13
6.3.1	Transformation, Inverse Transformation . . . . .	13
6.3.2	Implementationsüberlegungen . . . . .	13
<b>7</b>	<b>Burrows-Wheeler transform</b>	<b>14</b>
7.1	Definitionen . . . . .	14
7.2	Erläuterung . . . . .	14
7.3	Inverse Transformation . . . . .	14
7.4	Implementierung . . . . .	15
7.4.1	Transformation . . . . .	15
7.4.2	Invers . . . . .	15
7.4.3	Implementationsüberlegungen . . . . .	15

## Abbildungsverzeichnis

1	Verlustbehaftete Kompression eines Bildes . . . . .	5
2	Counter -und Zeichenbytes beim Run-Length . . . . .	6

## 1 Was ist Bzip2?

- Verlustfreies Datenkompressionsverfahren
- Ähnlich zu ZIP
- Basiert auf Burrows-Wheeler transform
- Autor: Julian Seward
- Veröffentlichung: 18. Juli 1996
- Verwendete Algorithmen
  - Run Length Encoding
  - Huffman Encoding
  - Move-to-Front Transform
  - Burrows Wheeler Transform

## 2 Was ist eine Datenkompression?

Als Datenkompression wird ein Vorgang bezeichnet, bei dem Daten so verarbeitet werden, dass sie mit weniger Bits darstellbar sind. Dabei kann die Kompression entweder **verlustfrei** oder **verlustbehaftet** sein.

### 2.1 Verlustfreie und verlustbehaftete Kompressionen

Für eine Kompression  $k : E \rightarrow A$ , die eine Eingabe  $E \in \{0, 1\}^*$  in eine Ausgabe  $A \in \{0, 1\}^*$  komprimiert, gilt:

$$k \text{ ist verlustfrei} \Leftrightarrow k \text{ ist invertierbar.}$$

Bei **verlustfreien** Kompressionen wird also die Forderung gesetzt, dass die komprimierte Ausgabe wieder in die originale Eingabe dekomprimiert werden kann.

Dass eine Kompression verlustbehaftet ist, ist nicht immer schlimm. Wichtig ist hierbei nur, dass die verlorene Information für den Menschen nicht erkennbar oder zumindest vernachlässigbar ist. So kann zum Beispiel bei einem Bild das

Farbspektrum verkleinert werden, ohne dass das menschliche Auge einen Unterschied bemerkt.



Abbildung 1: Ein Bild, welches mit dem mit dem verlustbehafteten JPEG-Verfahren komprimiert wurde (rechts), sodass Kompressionsartefakte zu erkennen sind [1].

In Texten hingegen können schon kleine Informationsverluste dafür sorgen, dass diese nicht mehr lesbar sind.

## 2.2 Entropiekodierungen und Substitutionskompressionen

Bei verlustfreien Kompressionen unterscheidet man grundsätzlich zwischen den **Entropiekodierungen** und **Substitutionskompressionen**.

**Substitutionskompressionen** versuchen, oft wiederholende Zeichenfolgen zusammenzufassen oder Muster und somit Redundanz in der Eingabe zu eliminieren. Bekannte Vertreter sind hier das LZ77 Verfahren und das im folgenden Abschnitt analysierte Run-Length Verfahren.

**Entropiekodierungen** arbeiten im Gegensatz zu Substitutionskompressionen nicht mit Zeichenfolgen oder Mustern in der Eingabe, sondern mit den relativen Häufigkeiten der einzelnen Zeichen. Dabei ist das Ziel, dass häufig auftretende Zeichen weniger Bits zur Folge haben und seltene Zeichen mehr Bits. Die Huffman- und arithmetische Kodierung sind Beispiele von Entropiekodierungen.

## 3 Run-Length encoding

Die Run-Length Kodierung ist eine Substitutionskompression, bei der sich wiederholende Zeichenfolgen zusammengefasst werden.

### 3.1 Kodierung

Sei  $aw \in \Sigma^*$  die Eingabe, wobei  $a = x^n$ ,  $x \in \Sigma$  eine sich wiederholende Zeichenfolge der Länge  $n$  ist. Das Teilwort  $a$  wird zu  $nx$  umgeformt, wonach der Algorithmus mit  $w$  fortführt. Folgendes ist ein Beispiel der Kodierung:

$$\text{jjjjjjkkijjj} \Rightarrow 7j\text{kkijjj} \Rightarrow 7j2kijjj \Rightarrow 7j2k1ijjj \Rightarrow 7j2k1i3j$$

### 3.2 Dekodierung

Sei  $n x w$  mit  $n \in \mathbb{N}$ ,  $x \in \Sigma$  und  $w \in (\mathbb{N} \cup \Sigma)^*$  die Eingabe. Das Teilwort  $n x$  wird zu  $x^n$  umgeformt, wonach der Algorithmus mit  $w$  fortführt. Folgendes ist ein Beispiel der Kodierung:

$2j1k6i3k \Rightarrow jj1k6i3k \Rightarrow jjk6i3k \Rightarrow jjkiiii3k \Rightarrow jjkiiiiikk$

### 3.3 Effizienz

Solange die Eingabe viele sich wiederholende Zeichenfolgen enthält, ist die Run-Length Kodierung sehr effektiv. Falls aber nicht, kann die Kompression sogar zu einem größeren Platzverbrauch führen. Die Zeichenfolge **abab...** würde zum Beispiel zu **1a1b1a1b...** komprimiert werden. Falls sowohl Counter als auch Zeichen genau ein Byte Platz bräuchten (kommt ganz auf die Implementierung an), würde die Kompression zu einer Verdopplung der Bits führen.

### 3.4 Implementierung

Eine einfache und effiziente Implementierung von Run-Length arbeitet mit 8-Bit Zeichen, meist ASCII, und hat für den Counter eine Obergrenze von 255, also auch 8-Bit. Das vereinfacht das Dekodieren deutlich, da Counter von Zeichen anhand der Bitposition des Textes unterschieden werden kann. Wenn das Byte auf einem ungeraden Index ist, handelt es sich bei diesem um einen Counterbyte und ansonsten um ein Zeichenbyte.

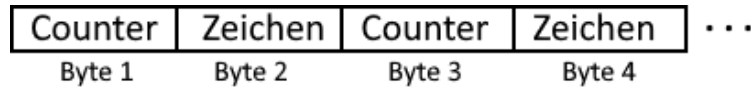


Abbildung 2: Counter -und Zeichenbytes beim Run-Length nach obiger Variante

Außerdem gibt es eine Variante von Run-Length, die zum Beispiel nur Zeichenwiederholungen der Länge 4 oder größer zusammenfasst, um die mögliche Verdopplung des Platzverbrauchs zu verhindern. **abab...** würde dann nicht verändert werden, da alle Zeichenwiederholungen die Länge 1 haben. So ist es aber nicht mehr einfach möglich, Counter von Zeichen zu unterscheiden.

## 3.5 Pseudocode

### 3.5.1 Kodierung

```
function run_length(input):  
    out  $\leftarrow$   $\epsilon$   
    counter  $\leftarrow$  1  
    run_char  $\leftarrow$  input[1]  
    for i in 2..n:  
        if input[i-1] = input[i]:  
            counter++  
        else:  
            out.append(counter, run_char)  
            counter  $\leftarrow$  1  
            run_char  $\leftarrow$  input[i]  
    out.append(counter)  
    out.append(run_char)  
    return out
```

### 3.5.2 Dekodierung

```
function run_length_decode(input):  
    out  $\leftarrow$   $\epsilon$   
    for (count, char) in input:  
        out.append(char * count)  
    return out
```

## 4 Huffman encoding

Die Huffman-Kodierung errechnet für Zeichen  $z \in \Sigma$  einer Eingabe einen Code

$$c : \Sigma \rightarrow \bigcup_{n \geq 1} \{0, 1\}^n.$$

Dafür zählt die Kodierung die relative Häufigkeit jedes Zeichens und gibt Zeichen mit hoher relativer Häufigkeit kürzere Codes und Zeichen mit kleinerer relativer Häufigkeit entsprechend längere Codes. Bei Eingaben mit ungleichen Häufigkeitsverteilungen machen einige wenige Zeichen den Großteil der Eingabe aus, sodass ein Großteil der Zeichen kurze Codes hat, was eine kleine durchschnittliche Codelänge und somit große Platzeinsparungen zur Folge hat.

### 4.1 Präfixfreiheit

Ein Code  $c : \Sigma \rightarrow \Gamma$  gilt dann als **präfixfrei**, wenn

$$\forall x, y \in \text{Bild}(c) : \nexists z \in \Gamma : xz = y.$$

Präfixfreie Codes wie bei der Huffman-Kodierung haben den Vorteil, man so aus einer Codefolge direkt die originale Eingabe linear errechnen kann. Eine Bitfolge ist genau dann ein Zeichen der originalen Eingabe, wenn die Bitfolge im Bild von  $c$  ist, welche bei der Huffman-Kodierung aus der gespeicherten Tabelle berechnet werden können.

### 4.2 Kodierung

Die Kodierung besteht aus folgenden Schritten:

1. Zähl in der Eingabe, wie oft jedes Zeichen vorkommt und berechne so die relative Häufigkeit.
2. Erstelle einen Graph, wobei jedes Zeichen einen Knoten hat, indem seine relative Häufigkeit und sein Zeichen gespeichert sind.
3. Solange es mehr als einen Knoten gibt, der kein Elternknoten besitzt:
  - Erstelle einen Knoten  $k$ , dessen beiden Kinder die Knoten im Graph mit der kleinsten relativen Häufigkeit sind.
  - Die linke Kante wird mit 0 und die rechte mit 1 markiert.
  - Setze die relative Häufigkeit von  $k$  auf die Summe der relativen Häufigkeiten seiner beiden Kinder.
4. Erstelle eine Tabelle, in der jedes Zeichen einem Code zugeordnet wird.
5. Traversiere den Baum von der Wurzel
  - Sobald ein Blatt erreicht wird, wird dessen Zeichen der zurückgelegte Weg als Kodierung zugewiesen.
6. Setze jedes Zeichen der Eingabe auf seinen Code in der Tabelle.



### 4.3 Dekodierung

Da bei der Kodierung der Eingabe auch die Tabelle mitgespeichert wird, kann man einfach den kodierten Text Bit für Bit durchgehen, bis die Bitfolge als Eintrag in der Tabelle zu finden ist. Da der Code präfixfrei ist, kann es keinen anderen Code in der Tabelle geben, der mit dieser Bitfolge erreicht werden könnte, sodass die Zuordnung eindeutig ist.

### 4.4 Effizienz

Da die Tabelle gespeichert werden muss und seine Größe linear von der Alphabetgröße abhängt, weil jedes Zeichen im Alphabet ein Eintrag besitzt, hat die Huffman-Kodierung einen riesigen Overhead. Bei einem Alphabet mit 256 Zeichen lohnt sich die Huffman-Kodierung für eine Eingabe mit weniger als 1000 Zeichen überhaupt nicht, da die Tabellengröße einen großen Anteil der Gesamtgröße ausmacht. Erst bei einer sehr großen Anzahl an Zeichen, z.B. 100000 ist die Tabellengröße vernachlässigbar.

### 4.5 Pseudocode

Da im 3. Schritt der Kodierung immer die Kinder mit der kleinsten relativen Häufigkeit gesucht werden, eignet sich ein Min-Heap als Datenstruktur für eine schnelle Suche.

#### 4.5.1 Tabellengenerierung

```
function huffman_table(input):
    table: string  $\Rightarrow$  string
    heap  $\leftarrow$  MinHeap()
    for char in input:
        heap.insert(char, occurrence(char))
    root  $\leftarrow$  {}
    while heap.hasItem():
        (char1, char1_occ)  $\leftarrow$  heap.pop()
        if heap.empty(): break
        (char2, char2_occ)  $\leftarrow$  heap.pop()
        root  $\leftarrow$  (char1_occ + char2_occ)
        heap.insert(root)
    for leaf in root:
        table.map(leaf.path  $\rightarrow$  leaf.char)
    return table
```

#### 4.5.2 Kodierung

```
function huffman_encode(input):  
    table ← huffman_table(input)  
    out ←  $\epsilon$   
    for char in input:  
        out.append(table(char))  
    return table, out
```

#### 4.5.3 Dekodierung

```
function huffman_decode(table, input):  
    buffer ←  $\epsilon$   
    out ←  $\epsilon$   
    for char in input:  
        buffer.append(char)  
        if buffer in table:  
            out.append(table(buffer))  
            buffer ←  $\epsilon$   
    return out
```

## 5 Transformation

### 5.1 Was ist eine Transformation?

Transformationen sind Funktionen auf einer Menge von Eingabewörtern  $\Sigma^*$ , welche diese auf eine, oft gleiche, Menge von Ausgabewörtern  $\Sigma'^*$  abbildet. Die Eingabe und das Ergebnis sind dabei jedoch gleich lang.

Sei  $T$  eine Transformation mit

$$T : \Sigma^* \rightarrow \Sigma'^*, w \mapsto T(w) \quad (1)$$

Die Anwendung einer Transformation wird im folgenden für ein Wort  $w \in \Sigma^*$  gekürzt beschrieben.

$$w^T := T(w) \quad (2)$$

Aus der oben genannten Definition folgt

$$|w^T| = |w| \quad (3)$$

Beispiele simpler Transformationen:

$$\begin{aligned} abcd &\xrightarrow{R} dcba \quad (\text{Umkehrung}) \\ dbca &\xrightarrow{\text{sort}} abcd \quad (\text{Lexikographische Sortierung}) \\ abcd &\xrightarrow{0} 0000 \quad (\text{Nulling}) \end{aligned}$$

Man beobachte: die Permutation ist ein Spezialfall der Transformation.

### 5.2 Warum sind Transformationen nützlich?

Das Ziel der Datenkompression ist logischerweise das Einsparen von Speicherplatz. Wenn jedoch die Länge der Ausgabe einer Transformation der Länge der Eingabe gleicht, wozu benötigen wir sie dann?

Transformationen können benutzt werden, um Eingaben in eine für eine Kodierung (o.Ä.) günstige Form zu bringen. Dies steigert die Effizienz der Kodierung und spart somit mehr Speicherplatz.

Wir stellen die Forderung, dass eine Transformation invertierbar sein muss, um für (verlustfreie) Datenkompression nützlich zu sein. Die Umkehrfunktion einer Transformation  $T$  wird wie folgt definiert.

$$T^{-1} : \Sigma'^* \rightarrow \Sigma^*, w = (w^T)^{T^{-1}}$$

Jetzt gilt es nur noch, tatsächlich nützliche Transformationen zu finden und anzuwenden.

## 6 Move-to-front transform

### 6.1 Erläuterung

MTF ist eine Transformation, welche Zeichen der Eingabe auf deren Index in einer Liste "neulich genutzter" Zeichen abbildet.

Sei  $\Sigma$  das Eingabealphabet mit  $|\Sigma| = n$ . Wir initialisieren die Liste  $A(1..n)$  folgendermaßen.

$$A_{1..n} = (\sigma_1, \dots, \sigma_n), \sigma_1, \dots, \sigma_n \in \Sigma \text{ paarweise Verschieden} \quad (4)$$

$A$  enthält initial das gesamte Alphabet, per konvention lexikographisch sortiert.

Wir definieren die Hilfsfunktion  $t_{\text{MTF}}$ :

$$\begin{aligned} t_{\text{MTF}} : \Sigma^n \times \Sigma &\rightarrow \Sigma^n \times \Sigma, \\ t_{\text{MTF}}(A, \sigma) &:= (A_i \circ (A \setminus A_i), i), \quad i = \text{Index von } \sigma \text{ in } A \end{aligned} \quad (5)$$

$t_{\text{MTF}}$  beschreibt einen Schritt des MTF-Algorithmus. Für ein Zeichen  $\sigma \in \Sigma$  der Eingabe wird zunächst dessen Index in  $A$  in die Ausgabe geschrieben. Anschließend wird das Zeichen vorne an  $A$  gehängt.

Mithilfe von  $t_{\text{MTF}}$  lässt sich MTF nun definieren.

$$\begin{aligned} \text{MTF}_A : \Sigma^* &\rightarrow \Sigma^* \\ w^{\text{MTF}_A} &:= i \circ w' \\ \text{mit } (A', i) &= t_{\text{MTF}}(A, w_1) \\ w' &= (w_2, \dots, w_k)^{\text{MTF}_{A'}} \end{aligned} \quad (6)$$

## 6.2 Inverse Transformation

Die Invertierung der MTF ist trivial, da gilt

$$\text{MTF}_A^{-1} = \text{MTF}_A \quad (7)$$

## 6.3 Implementierung

### 6.3.1 Transformation, Inverse Transformation

```
function mtf(input):  
    A[|Σ|] ← Σ  
    out ← ε  
    for char in input:  
        char_index ← A.indexOf(char)  
        out.append(char_index)  
        delete A[char_index]  
        prepend(char, A)  
    return out
```

### 6.3.2 Implementationsüberlegungen

Da  $A$  für Transformation und Inversion gleich sein muss, ist es effizienter, ein festes  $A$  in der Spezifikation des Algorithmus festzulegen.

Die Geschwindigkeit der Transformation wird dramatisch verbessert, wenn für  $A$  eine Linked List verwendet wird.

## 7 Burrows-Wheeler transform

### 7.1 Definitionen

Sei  $w \in \Sigma^*$  mit  $|w| = n$ . Wir definieren die  $i$ -te Rotation von  $w$  für  $1 \leq i \leq n$  folgendermaßen.

$$w^{\lambda_i} := (w_i, \dots, w_n, w_1, \dots, w_{i-1}) \quad (8)$$

Die Menge aller Rotationen von  $w$  sei

$$w^\lambda := \{w^{\lambda_i} : 1 \leq i \leq n\} \subset S_n \quad (9)$$

Ferner sei die Rotation von  $w$  des Ranges  $k$

$$w^{\Lambda_k} := w^{\lambda_i}, \text{ sodass } |\{w^{\lambda_j} : w^{\lambda_j} \leq_l w^{\lambda_i}, 1 \leq i, j \leq n\}| = k \quad (10)$$

In einer lexikographisch sortierten Liste von Rotationen ist  $w^{\Lambda_k}$  der  $k$ -te Eintrag.

### 7.2 Erläuterung

BWT ist eine Transformation, welche die durchschnittliche Run-Length in der Eingabe erhöht. Dazu wird die lexikographische Rangfolge der Zeichen in betracht gezogen.

BWT sei für eine Eingabe  $w \in \Sigma^*$  definiert als

$$\begin{aligned} \text{BWT} : \Sigma^* &\rightarrow (\Sigma \cup \{\$\})^*, \\ w^{\text{BWT}} &:= (u_{n+1}^{\Lambda_1}, \dots, u_{n+1}^{\Lambda_{n+1}}) \\ &\text{mit } u = w \circ \$ \end{aligned} \quad (11)$$

Aus einer sortierten Liste von Rotationen der Eingabe wird die Ausgabe jeweils aus dem letzten Zeichen jeder Rotationen gebildet. Davor wird jedoch ein "\$" als Endzeichen an die Eingabe gehängt. Es gilt  $\$ \notin \Sigma$ .

### 7.3 Inverse Transformation

Für die Invertierung einer Ausgabe  $w \in \Sigma^*$  und dessen lexikographisch sortierte Permutation  $w'$  der Transformation wird zunächst eine Hilfspermutation definiert.

$$\gamma : \mathbb{N} \rightarrow \mathbb{N}, \quad w_{\gamma(x)} = w'_x \quad (12)$$

$\gamma$  bildet die Position eines Zeichens in  $w$  auf dessen Position in  $w'$  ab. Ein Zeichen wird also auf dessen neue Position abgebildet, nachdem  $w$  sortiert wurde.

Man beobachte:

$$\gamma \text{ besteht aus einem einzigen Zyklus der Länge } n + 1. \quad (13)$$

Nun lässt sich die Inversion von BWT folgendermaßen definieren.

$$w^{\text{BWT}^{-1}} := (w_{\gamma(\$)}, w_{\gamma(\gamma(\$)), \dots})^R \quad (14)$$

Merke: es gilt  $|w| = |w^{\text{BWT}^{-1}}|$

## 7.4 Implementierung

### 7.4.1 Transformation

```
function burrows_wheeler(input):
    rotations[n, n] ← input.rotations()
    rotations.sort()
    return rotations[1..n, n]
```

### 7.4.2 Invers

```
function inverse_burrows_wheeler(input):
    rotations[n, n]
    for k in n..1:
        rotations[1..n, k] ← input
        rotations.sort()
    return rotation ending with $
```

### 7.4.3 Implementationsüberlegungen

Suffix Arrays können in der Transformation in linearer Zeit konstruiert werden, um die Transformation letztendlich linear umzusetzen.

Die Platz-Effizienz der Pseudocode Beispiele ist relativ schlecht. Statt ganze Matrizen zu speichern, können Pointer in die jeweiligen Strings verwendet werden.

## Literatur

- [1] *JPEG Artefakte*. (besucht am 19.02.2023, 16:00). URL: [https://de.wikipedia.org/wiki/Datenkompression#/media/Datei:Jpegartefakt\\_jpegartefact.jpg](https://de.wikipedia.org/wiki/Datenkompression#/media/Datei:Jpegartefakt_jpegartefact.jpg).