



Wird dadurch Speicher eingespart?

Codierung Zeile1:	010000	neue Länge -10 Bit
Codierung Zeile2:	000111 100001 001000	neue Länge +2 Bit
Codierung Zeile3:	000110 100001 000001 100001 000111	neue Länge +14 Bit
<b>Total</b>		<b>neue Länge +6 Bit</b>

Total werden leider nach dieser "Kompression" 6 Bit mehr benötigt. Die Ursache liegt in der maximalen Anzahl von 16 Bit pro Zeile, welche nur einmal erreicht wird. Wenn wir nun stattdessen eine maximale Länge von 7 aufeinanderfolgenden Bit definieren erhalten wir folgendes Resultat:

Codierung Zeile1:	0111 0111 0010	neue Länge -4 Bit
Codierung Zeile2:	0111 1001 0111 0001	neue Länge -0 Bit
Codierung Zeile3:	0110 1001 0001 1001 0111	neue Länge +4 Bit
<b>Total</b>		<b>neue Länge -0 Bit</b>

Was lernen wir daraus?

- Da es im binären System kein Trennzeichen gibt, müssen wir eine fixe Länge für die Mengenangabe definieren. 1 Bit Daten, 5 Bit Mengenangabe bzw. 1 Bit Daten 3 Bit Mengenangabe in den obigen Beispielen
- Je nach Daten, welche wir komprimieren möchten schwankt das Resultat enorm! Zeile 1 und 2 sind eher komprimierbar, Zeile 3 kaum. -> Es gibt nicht komprimierbare Daten  
Wenn das Bild grösser wäre, würde sich das Komprimieren auf jeden Fall lohnen!
- Eine vorgängige Analyse der zu komprimierenden Daten kann bei der Auswahl des Verfahrens nützlich sein. (In diesem Fall um die Länge der Mengenangabe zu optimieren)
- Ein Komprimierungsverfahren hat immer einen Anteil **Daten** und einen Anteil **Zusatzinformation**, welche aus dem Verfahren heraus entsteht. Das Ziel ist es nun, die **Nutzdaten** dermassen zu verkleinern, dass zusammen mit den **Zusatzinformationen** am Ende ein kürzerer Code entsteht.

#### Laufängencodierung mit Textdaten

Bei Textdaten (erweiterter ASCII, ANSI usw.) kann die Laufängencodierung ebenfalls eingesetzt werden. Um Speicherplatz zu sparen wird die Mengenangabe binär codiert. Selten kommen mehr als 2 identische Zeichen hintereinander vor im normalen Sprachgebrauch.

Idee: 1 Zeichen + binäre Mengenangabe (0 -> kommt 1x vor, 1 --> kommt 2x vor)

Beispiel:

100 leere Brunnen	$1_2 0_2 1_2 0_2 1_2 0_2 e_2 1_2 r_2 0_2 e_2 0_2 0_2 B_2 0_2 r_2 0_2 u_2 0_2 n_2 1_2 e_2 0_2 n_2 0_2$
17x8 Bit = 136 Bit	$(14 \times 8) + (14) \text{ Bit} = 112 + 14 \text{ Bit} = 126 \text{ Bit}$

Je nach Text kann die Laufängencodierung eine Ersparnis bringen. Als Beispiele finden man deshalb oft solche "Texte": XXXXX YYYYY ZZZZZZZZ

Zum Thema Laufängencodierung gibt es natürlich Übungen →1

#### Wörterbuchverfahren zur Datenkompression S. 184

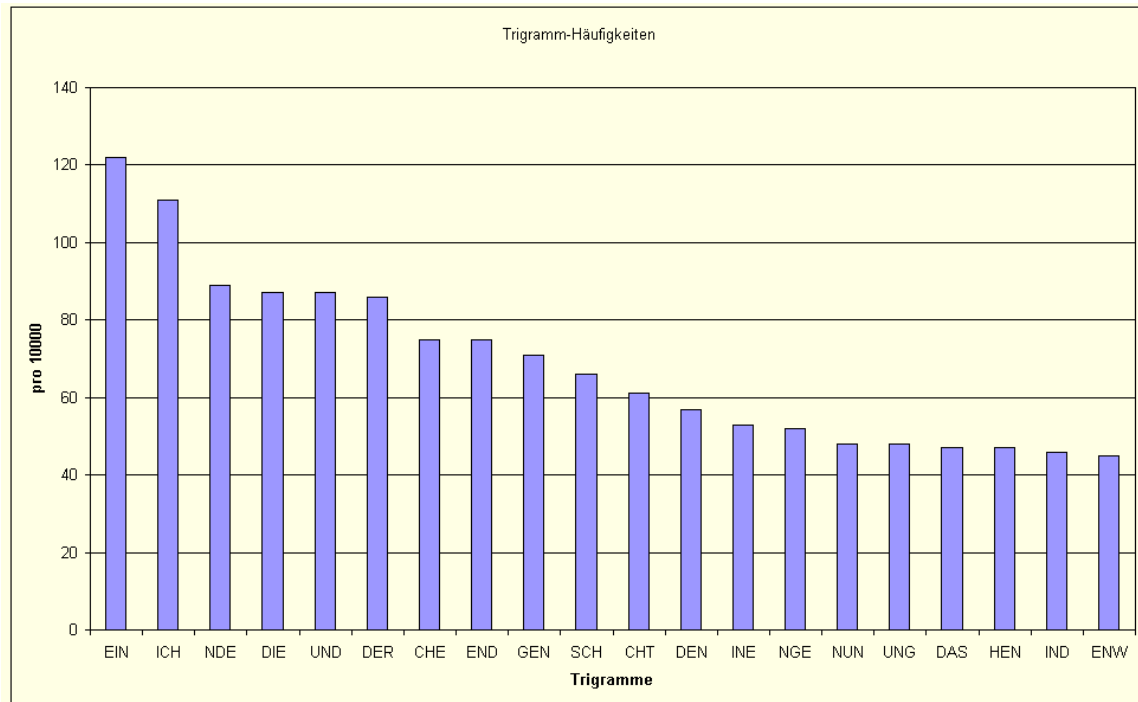
Bei der Wörterbuchmethode (engl. token based compression) werden wiederkehrende, häufig auftretende Muster zusammengefasst und durch einen kürzeren Code **ersetzt**. Dazu wird eine Liste (Wörterbuch) benötigt, in welcher die Codes und ihre ursprüngliche Bedeutung abgelegt sind. Übertragen wird dann Das **Wörterbuch** und die **codierten Daten**.

Die minimalste Variante eines Wörterbuchs für Text könnte so aussehen:

<Leer>: 00000      A: 00001      B: 00010      C: 00011 .....      Z: 11010

Mit dieser Methode kann ein 8 Bit ASCII Zeichen auf lediglich 5 Bits dargestellt werden, zudem müsste das Wörterbuch nicht einmal mitgeliefert werden, da es bei der Dekompression als bekannt angenommen werden könnte.

Ein Wörterbuch kann statt einzelner Buchstaben auch ganze Wörter aufnehmen. So kann beispielsweise für die deutsche Sprache ein intelligentes Wörterbuch erstellt werden anhand der statistischen Vorkommen von Buchstabenkombinationen:



(Quelle: [http://www.staff.uni-mainz.de/pommeren/Kryptologie/Klassisch/1\\_Monoalph/deutsch.html](http://www.staff.uni-mainz.de/pommeren/Kryptologie/Klassisch/1_Monoalph/deutsch.html))

Anhand dieser Tabelle könnte das Wörterbuch von oben um die Einträge aus dem Diagramm erweitert werden.

Bsp: EIN: 11011 (27) ICH: 11100 (28) NDE: 11101 DIE: 11110 UND: 11111 (31)

Der Satz "ICH UND MEIN HUND " würde dann so dargestellt: (28)(0)(31)(0)(13)(27)(0)(8)(31)  
Dadurch würden lediglich  $9 \times 5 = 45$  Bit benötigt, im Original wären es mit ASCII  $17 \times 8 = 136$  Bit.

Was passiert aber mit englischen Texten oder mit neutralen Daten? (Bilder, Programme usw.)  
Dazu gibt es spezialisierte Algorithmen (Rechenmodelle), welche vor der Kompression die Daten analysieren und dann selber eine Häufigkeit von wiederholenden Mustern erstellen. (siehe weitere Themen für Interessierte)

### Huffman Code (David A. Huffman 1952) S. 184

Der Huffman Code funktioniert auch nach dem System der Zeichenersetzung. Dabei werden die zu komprimierenden Daten analysiert, um herauszufinden welches Zeichen wie häufig vorkommt. Im Gegensatz zu den bisher betrachteten Verfahren ist beim Huffman Code keine fixe Länge der codierten Zeichen notwendig. Anhand des Aufbaus des Codes kann bei **sequentiell**em Durchlaufen entschieden werden wann ein neues Zeichen kommt.

Häufige Zeichen erhalten einen kurzen Code, seltene Zeichen einen längeren Code.  
Als Hilfsmittel zur Erstellung der Codes wird ein Baum erstellt (Metadaten). Am einfachsten ist das Verfahren an einem Beispiel zu verstehen:

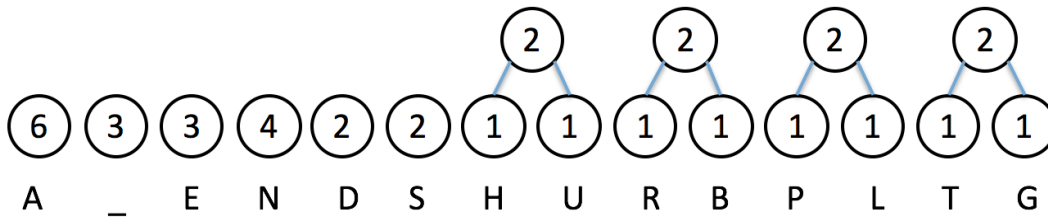
Zu komprimierende Daten: *DAS HAUS DER BANANENPLANTAGE*

Zählen der Vorkommen: (Das Leerzeichen nicht vergessen, wir schreiben es als "\_")

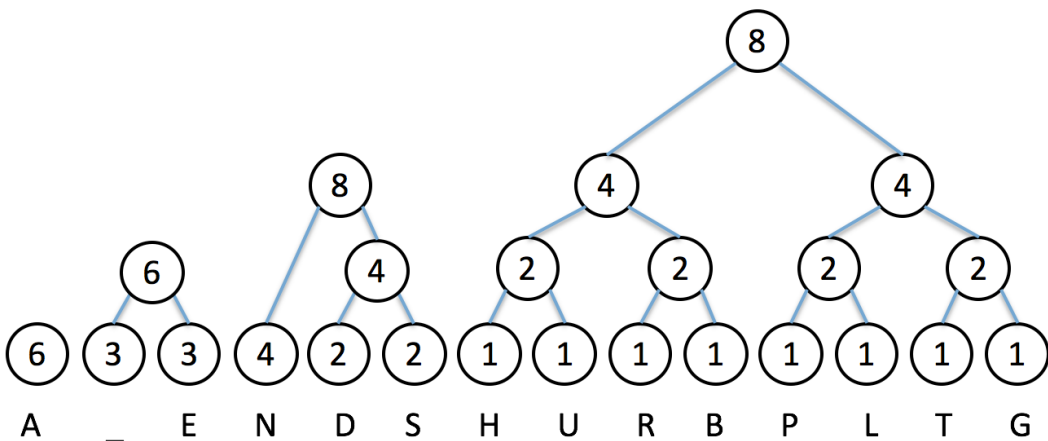
D: 2 A: 6 S: 2 \_: 3 H: 1 U: 1 E: 3 R: 1 B: 1 N: 4 P: 1 L: 1 T: 1  
G: 1

Tip: Am besten addieren wir alle Vorkommen und vergleichen Sie mit den Anzahl Zeichen der Daten!

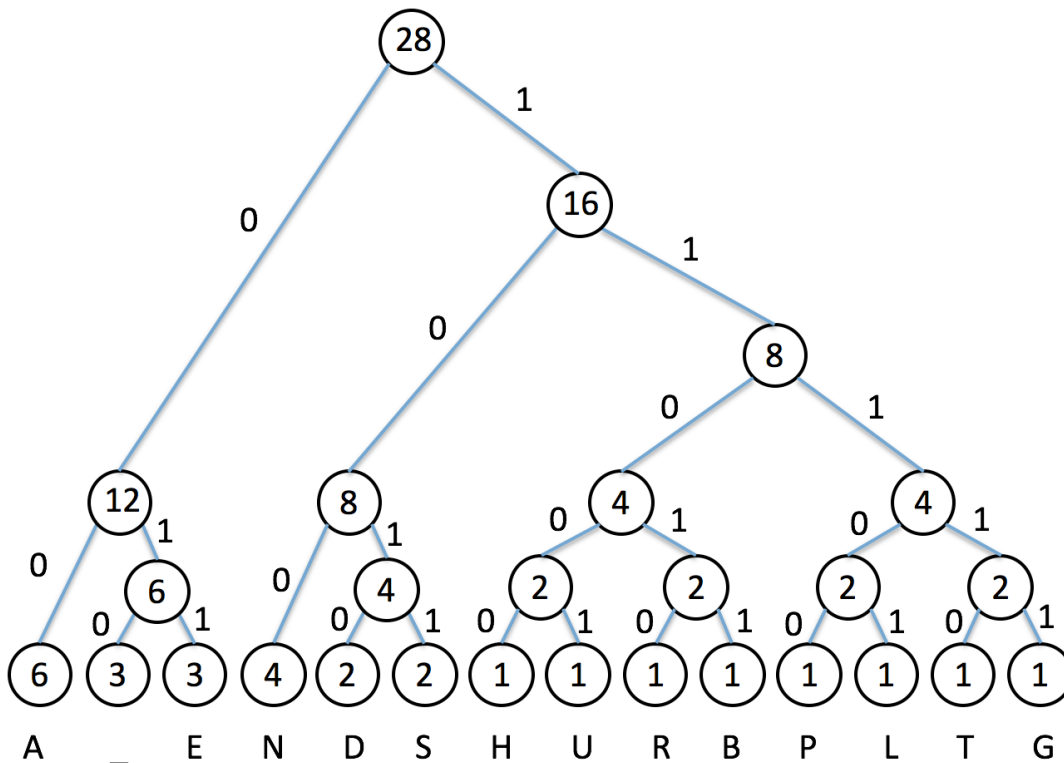
Nun bauen wir Schritt für Schritt den Baum auf, indem wir immer die tiefsten Vorkommen zusammenziehen zu einem neuen Knoten im Baum. Der neue Knoten erhält die Summe seiner beiden Äste als Inhalt:



Die Äste werden immer weiter verknüpft. Wichtig: immer die tiefsten Nummern miteinander verbinden!



Wenn der Baum bis oben zur Wurzel fertiggestellt ist, werden die Äste links und rechts jeweils mit 0 oder 1 beschriftet:



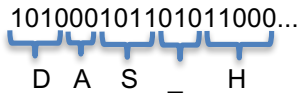
An dieser Stelle wäre wieder einmal eine Kontrolle des Totals (28) fällig.

Der Baum dient nun dazu die ursprüngliche Nachricht zu komprimieren:

**DAS HAUS DER BANANENPLANTAGE**

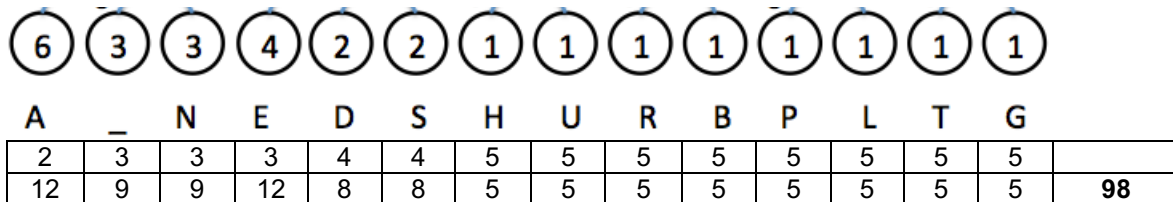
Die einzelnen Zeichen können nun vom Baum abgelesen werden und bilden zusammengesetzt den komprimierten Code.

D	1010
A	00
S	1011
_	010
H	11000
...	



Das Ziel, häufig auftretenden Zeichen einen kurzen Code zu vergeben wurde erreicht. Das "A", welches am häufigsten vorkommt hat den kürzesten Code "00" erhalten.

Wie viel Speicher wird nun benötigt? Dazu können wir nochmals den Baum zur Hilfe nehmen. Wir notieren unter jedes Zeichen die Länge des Codes (Aster) und multiplizieren diese mit der Anzahl Vorkommen. (diese haben wir ja schon)



Es werden also 98 Bits verwendet, diese zählen zu den **Nutzdaten**. **Der Baum** (Metadaten bzw. das Wörterbuch) müssten zusätzlich noch übertragen werden. Wenn wir nun nur die Nutzdaten anschauen, sparen wir mit 98 Bit gegenüber dem erweiterten ASCII (28 \* 8 Bit = 224 Bit) enorm Speicherplatz ein.

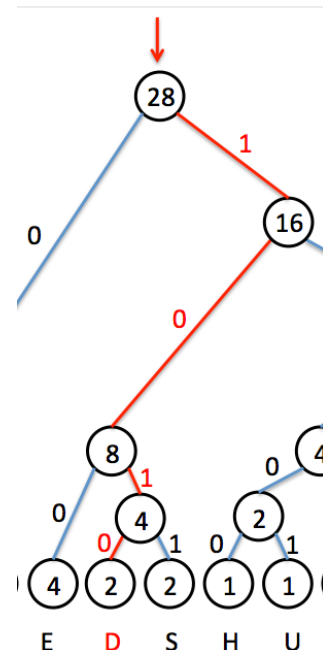
**Dekompression des Huffman Codes**

Bei Dekomprimieren von Inhalten, welche nach dem Huffman Verfahren verkleinert wurden benötigt man 2 Angaben: Die **Nutzdaten** (Bitfolge) Den **Baum**

Nutzdaten:  
1010001011010110000011001101101010101001101001011011 ...

Um die Originaldaten herzustellen nehmen wir die Bits aus den Nutzdaten und folgen dem Pfad, welcher uns die Nutzdaten vorgeben. Es wird immer oben im Baum begonnen:

Rechts ist sichtbar wie die Wiederherstellung funktioniert. Der Code der Nutzdaten "1010" wird durchlaufen bis ein Blatt (Buchstabe) erreicht ist. In diesem Fall der Buchstabe "D". Nach diesen 4 Bit wird wieder oben im Baum angefangen usw. Dies geschieht solange bis alle Nutzdaten "ausgepackt" worden sind. Wir erkennen, dass die Dekodierung nur möglich ist, wenn am Anfang der Nutzdaten begonnen wird! ➔





Auf linuxbasierten Systemen wird anders komprimiert, da finden wir .tar.gz Dateien. Auch dieses (zusammengesetzte) Format übernimmt die beiden Funktionen:

1. tar: Zusammenfassen mehrerer Dateien und Verzeichnisse in eine Datei (Archiv)
2. gz: Reduzieren des Speicherplatzes durch Kompression

tar ist die Abkürzung für "Tape Archive", welches für die alten Magnetbänder steht. Anders als beim .zip Format werden im .tar die einzelnen Dateien direkt hintereinandergereiht, ohne eine zusätzliche Ordnerstruktur innerhalb des Archivs abzubilden. Selbstverständlich können aber auch Verzeichnisse und Dateien in einem Archiv enthalten sein, lediglich die interne Speicherung ist unterschiedlich. Das gz Format zum Komprimieren funktioniert praktisch gleich wie bei .zip, nämlich mit dem Deflate Algorithmus.

Auch unter Linux kann in der grafischen Oberfläche einfach mit .tar.gz Dateien umgegangen werden. Da es aber sehr peinlich werden kann, wenn man in der Shell nicht mit diesen Dateien umgehen kann, üben wir dies. →4  
Hilfreich ist diese Tabelle:

Optionen von tar	
Option	Beschreibung
--help	Zeigt eine vollständige Übersicht über alle Optionen.
--version	Gibt die installierte Version von tar aus.
-c	Ein neues Archiv erzeugen.
-d	Dateien im Archiv und im Dateisystem miteinander vergleichen.
-f	Archiv in angegebene Datei schreiben. / Daten aus angegebener Datei lesen.
-j	Archiv zusätzlich mit <b>bzip2</b> (de)komprimieren.
-J	Archiv zusätzlich mit <b>xz</b> (de)komprimieren.
-k	Das Überschreiben existierender Dateien beim Extrahieren aus einem Archiv verhindern.
-p	Zugriffsrechte beim Extrahieren erhalten.
-r	Dateien an ein bestehendes Archiv anhängen.
-t	Inhalt eines Archivs anzeigen.
-u	Nur Dateien anhängen, die jünger sind als ihre Archiv-Version.
-v	Ausführliche Ausgabe aktivieren. Hierbei ist zu beachten, dass man dies möglichst an den Anfang des Befehls anhängt, wenn mehrere Optionen kombiniert werden. Z.B. würde -cfv zu einer Fehlermeldung führen. Korrekt wäre -vcf.
-w	Jede Aktion bestätigen.
-x	Dateien aus einem Archiv extrahieren.
-z	Archiv zusätzlich mit <b>gzip</b> (de)komprimieren.
-Z	Archiv zusätzlich mit <b>compress</b> (de)komprimieren.
-A	Inhalt eines bestehenden Archivs in ein anderes Archiv kopieren.
-M	Mehrteiliges Archiv anlegen/anzeigen/extrahieren.
-L	Medium wechseln, wenn ZAHL KBytes geschrieben sind.
-W	Archiv nach dem Schreiben prüfen.

(Quelle: <http://wiki.ubuntuusers.de/tar>)

### Weitere Themen für Interessierte

Falls Ihnen das Thema zusagt, können Sie eigene Recherchen zu den folgenden Themen durchführen:

- LZW Algorithmus (siehe Video im Ordner)
- 4 Komprimierungs- und Archivlösungen im Vergleich (<http://www.tomshardware.de/archivierung-zip-kompression,testberichte-240497.html>)







