

# ADS: Algorithmen und Datenstrukturen 2

## Teil 1

Prof. Dr. Gerhard Heyer

Institut für Informatik  
Abteilung Automatische Sprachverarbeitung  
**Universität Leipzig**

11. April 2018

[Letzte Aktualisierung: 11/04/2018, 11:17]

Übungsaufgaben, Vorlesungsfolien, Termin- und Raumänderungen, ...

<http://adsprak.informatik.uni-leipzig.de>

Erreichbar über

<http://asv.informatik.uni-leipzig.de/ss18/>

→ Lehre → Algorithmen und Datenstrukturen 2

## Anmeldung für Vorlesungen und Übungen

### Bitte befolgen Sie die Regeln

- ➊ Melden Sie sich zur **Vorlesung** im **AlmaWeb** an.
- ➋ Melden Sie sich zu einer der **Übungen** im **AlmaWeb** an.
- ➌ Beachten Sie die Anmeldefristen!

### Achtung!

- Ohne Anmeldung haben sie **keinen** Platz in einer der Übungen.
- Falls Sie sich **nicht rechtzeitig im AlmaWeb** angemeldet haben, so ist das Studienbüro ihr Ansprechpartner, **nicht** wir.
- Wechsel in andere Übungsgruppe nur mit Tauschpartner – um den **Sie** sich kümmern müssen.

- Abzugeben sind Lösungen zu fünf Aufgabenblättern.
- Termine

Aufgabenblatt	1	2	3	4	5+6
Ausgabe	18.04.	02.05.	16.05.	30.05.	13.06.
Abgabe	25.04.	09.05.	23.05.	30.05.	20.06.

- Lösungen sind **spätestens direkt vor** Beginn der Vorlesung im Hörsaal abzugeben. Spätere Abgaben werden **nicht** gewertet
- Lösungen werden bewertet und in der auf den Abgabetermin folgenden Übungsstunde zurückgegeben.

## Vorläufige Termine der Übungsgruppen

Tag	Uhrzeit	Raum	Übungsleiter
Di	09:15 - 10:45	SG 3-13	Sven Findeiß
Di	11:15 - 12:45	SG 3-13	Sven Findeiß
Di	11:15 - 12:45	SG 3-11	Martin Reckziegel
Mi	11:15 - 12:45	SG 3-13	Thomas Efer
Mi	11:15 - 12:45	Härtelstraße 16-18, S 109	Manuela Geiß
Mi	13:15 - 14:45	SG 3-13	Thomas Efer
Mi	13:15 - 14:45	Härtelstraße 16-18, S 109	Manuela Geiß
Mi	15:15 - 16:45	Härtelstraße 16-18, S 109	Manuela Geiß
Mi	17:15 - 18:45	Härtelstraße 16-18, S 109	Manuela Geiß
Do	13:15 - 14:45	SG 3-11	Martin Reckziegel
Do	13:15 - 14:45	SG 3-13	Christian Kahmann

## Zulassungsvoraussetzungen

- rechtzeitige(!) Modul- UND Prüfungsanmeldung in **AlmaWeb**
- Erreichen von 50% der Punkte in den Übungsaufgaben

- Team-Arbeit bei den Übungsblättern ist prinzipiell gestattet, ABER:
- jeder sollte alle abgegeben Lösungen erklären können;
- jeder muß eine eigene Abgabe haben;
- markieren Sie Gruppenarbeiten;
- **keine Kopien** von Ausarbeitungen;
- Abgabe in Papierform, **keine Möglichkeit** per email abzugeben;
- **Lösungen links oben tackern**; kein: kleben, vernähen, Büroklammern, schweißen, sonstwas;
- Lose Blätter können Punkteverlust bedeuten;
- **Name, Matrikelnummer und Übungsgruppe** auf jedes Blatt.

# Abgabe vor der Vorlesung, sortiert nach Übungsgruppen

Wenn Sie direkt vor der Vorlesung nicht abgeben können:

- geben Sie die Abgabe einem Ihrer Kommilitonen mit zur Abgabe direkt vor der Vorlesung
- oder werfen Sie die Abgabe bis spätestens **16 Uhr** (sechzehn Uhr) am Vortag in den Briefkasten der Abteilung *Automatische Sprachverarbeitung* in der 5. Etage im Augusteum, Raum A514.

## Spätere Abgaben werden nicht gewertet!



<http://adsprak.informatik.uni-leipzig.de/>

Kontaktadresse

[adshelp@informatik.uni-leipzig.de](mailto:adshelp@informatik.uni-leipzig.de)

Einführung in allgemeine *Entwurfsprinzipien* für Algorithmen, die sich für viele Problemstellungen anwenden lassen.

- Greedy-Algorithmen
- Dynamische Programmierung
- Probabilistische und heuristische Algorithmen (für NP-schwierige Probleme)

Anwendungsschwerpunkte

- Datenkompression
- Graphen und Netzwerke
- Rucksackproben, Travelling Salesman, ...

Bisheriges Hauptziel bei der Formulierung von Algorithmen

- schnelle Verarbeitung von Daten
- Effizienz bezüglich *Zeit*-Komplexität

Jetzt:

- *Speicher-effiziente Kodierung* (Darstellung)
- Datenkompression

Zwei grundsätzlich unterschiedliche Typen der Kompression

- **verlustfrei** (lossless)  
Codierung eindeutig umkehrbar.
- **verlustbehaftet** (lossy)  
Originaldaten können i.A. nicht eindeutig aus der komprimierten Codierung zurückgewonnen werden.  
Beispiele: JPEG, MPEG (mp3)

- **Laufängenkodierung:** Aufeinanderfolgende identische Zeichen werden zusammengefasst.
- **Huffman-Kodierung:** Häufiger auftretende Zeichen werden mit weniger Bits codiert.
- **LZW, LZ77, gzip:** Mehrfach auftretende Zeichenfolgen werden indiziert (*Wörterbuch/Codetabelle*), um dann durch ein einzelnes Zeichen dargestellt zu werden.

## Allgemeines Prinzip:

Ausnutzen von **Redundanz**

(Abweichungen von zufälliger Zeichenfolge, wiederkehrende Muster)

- Einfachster Typ von Redundanz:  
*Läufe* (runs) = lange Folgen sich wiederholender Zeichen.
- Beispiel:  
AAAABBBBAABBBBBCCCCCCCCDABCBAAABBBBCCCD.
- Ersetze jeden Run durch seine Länge und das wiederholte Zeichen:  
4A3B2A5B8CDABCB3A4B3CD
- Lohnt sich nur für Läufe mit Länge  $> 2$ .

**Wie können wir Text-Zeichen von Längenangaben unterscheiden?  
Insbesondere, wenn alle Zeichen (auch Ziffern) im zu codierenden  
Text vorkommen dürfen.**

# Laufängencodierung: Codierung mit Escape-Zeichen

- Wähle als **Escape-Zeichen** ein Zeichen Q, das in der Eingabe wahrscheinlich nur selten auftritt.
- Jedes Auftreten dieses Zeichens besagt, dass die folgenden beiden Buchstaben ein Paar (Zähler, Zeichen) bilden.
- Ein solches Tripel wird als *Escape-Sequenz* bezeichnet.
- Zählerwert  $i$  wird durch  $i$ -tes Zeichen des Alphabets dargestellt.
- Codierung erst für Läufe ab Länge 4 sinnvoll.
- Auftreten des Escape-Zeichens Q selbst muss speziell codiert werden, z.B. als Run der Länge 1 oder  $Q_{\_}$  ( $\_$ =Leerzeichen), etc.
- Beispiel:

Text	AAAABBBBAABBBBBCCCCCCCCDABCQC.
Codierung	QDABBBAAQEBQHCDABCQ_ <sub> </sub> C

Für Sequenzen von Bitwerten  $\in \{0,1\}$ :

- Beobachtung: Läufe von 0 und 1 wechseln sich ab.
- Nutze das aus: gebe nur noch Lauflängen an.
- Beispiel:

Sequenz	0001110111101100011111
Kodierung	3 3 1 4 1 2 3 5

**Was ist, wenn die Sequenz mit 1 beginnt?**



- Ansatzpunkt für Kompression: In Sprache kommen Zeichen mit unterschiedlicher Häufigkeit vor, z.B. **E** häufiger als **Y**.
- Beispiel: **ABRACADABRA**
- Standardcodierung (unkomprimiert) mit 5 Bits pro Zeichen:  
**00001 00010 10010 00001 00011 00001 00100 00001 00010 10010 00001**
- Häufigkeit der Buchstaben: 

A	B	C	D	R
5	2	1	1	2
- **Idee:** verwende für häufige Zeichen kurze Bitsequenzen, für seltene dafür längere. Minimiere so die Gesamtzahl der benötigten Bits.

**Wie kann man, trotz unterschiedlich langer Bitfolgen, erkennen welche Bits welches Zeichen codieren?**

**Gegeben sei folgende Codierung mit 0 und 1 (bzw. kurz oder lang wie im Morsealphabet):**

$$E = 0$$

$$T = 1$$

$$A = 01$$

Was bedeutet 0101?

Buchstabenkombination    ETA    AA    ETET    AET

**Ohne Sonderzeichen keine eindeutige Zerlegung!**

## Präfixcode (Präfix-freier Code)

- Kein Codewort ist Präfix eines anderen Codeworts.
- Durch Präfix-Freiheit: Codierung eindeutig umkehrbar.
- Im Beispiel (ABRACADABRA):

Buchstabe	A	B	C	D	R
Häufigkeit	5	2	1	1	2
Code	0	100	1010	1011	11

Codierung:

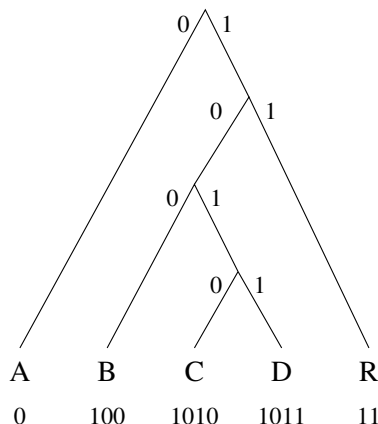
01001101010010110100110

AB R AC AD AB R A

# Codierung mit variabler Länge: Tries

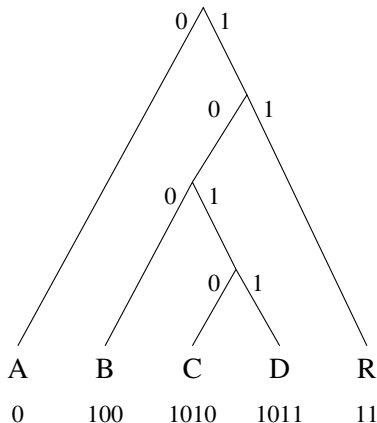
Binärer Präfixcode lässt sich durch Binärbaum (Trie) darstellen:

- **Blätter** = zu codierende Zeichen
- Zur **Codierung eines Zeichens**:  
laufe von Wurzel zum Blatt des Zeichens und gebe Kantenlabel aus (0 bei Verzweigung nach links; 1, nach rechts)
- Übliche Bezeichnung: **Trie**  
(dazu später mehr)



Gegeben: Bitsequenz und Präfixcode als Trie

- ① Beginne bei Wurzel im Baum
- ② Wiederhole bis zum Ende der Bitsequenz
  - ① Lies nächstes Bit. Bei 0 verzweige nach links im Baum, sonst nach rechts.
  - ② Falls erreichter Knoten ein Blatt ist, gib das Zeichen des Blatts aus und springe zurück zur Wurzel.



Der *Huffman-Code* ist ein präfix-freier Code mit variabler Länge, der die unterschiedlichen Zeichen-Häufigkeiten ausnutzt und systematisch erzeugt werden kann.

## Erzeugung des Huffman-Codes

**Schritt 1:** Zähle Häufigkeit der Zeichen in der zu codierenden Zeichenfolge

Das folgende Programm ermittelt die Buchstaben-Häufigkeiten einer Zeichenfolge `a` und trägt diese in ein Feld `count` ein.

```
for ( i = 0 ; i <= 26 ; i++ ) count[i] = 0;  
for ( i = 0 ; i < a.length ; i++ ) count[index(a[i])]++;
```

Dabei liefert `index(c)` den Index eines Zeichens `c`:  
`index(A)=1, index(B)=2, ...; index(_):=0`

“A SIMPLE STRING TO BE ENCODED USING  
A MINIMAL NUMBER OF BITS”

Dazugehörige Häufigkeits-Tabelle

	_	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
k	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26
count[k]	11	3	3	1	2	5	1	2	0	6	0	0	2	4	5	3	1	0	2	4	3	2	0	0	0	0	0

## **Schritt 2:** Aufbau des Tries (entsprechend den Häufigkeiten)

Für jedes Zeichen mit Häufigkeit  $\neq 0$ , erzeuge einen Baum jeweils mit dem Zeichen als einzigem Knoten.

## **Schritt 3:** Kombiniere iterativ die Bäume mit kleinsten Häufigkeiten

- Wähle zwei Bäume mit minimalen Häufigkeiten (Die Häufigkeit eines Baums ist die Summe der Häufigkeiten aller seiner Zeichen. Während der Erzeugung speichern wir Häufigkeit jedes Baums ab. )
- Verwende hierfür als Zwischenstruktur einen *MinHeap*
- Erzeuge neuen Knoten, der diese beiden Bäume als Nachfolger hat
- Iteriere Schritt 3 bis alle Knoten miteinander zu einem einzigen Baum verbunden sind.

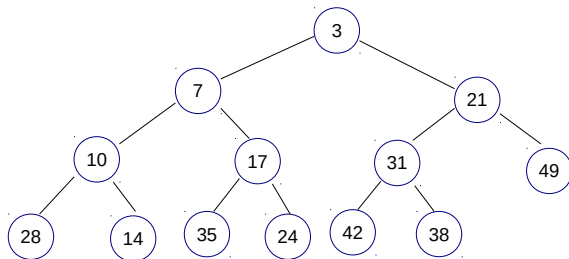
Am Ende sind Zeichen mit geringen Häufigkeiten weit unten im Baum; Zeichen mit großen Häufigkeiten nah an der Wurzel.



- Datenstruktur Heap („Halde“) als Spezialfall eines binären Baumes
- **Ein Heap ist ein Binärbaum**, der
  - die Heapeigenschaft hat (Kinder sind größer/kleiner als der Vater)
  - bis auf die letzte Ebene vollständig besetzt ist
  - höchstens eine Lücke in der letzten Ebene hat, die ganz rechts liegt
- Heaps i. A. genutzt, um den Datentyp *Vorrangwarteschlange* (engl. „priority queue“) zu implementieren
- Folgende Operationen sollen dabei unterstützt werden:
  - Lesen des Objektes mit der kleinsten/größten Priorität
  - Löschen des Objektes mit der kleinsten/größten Priorität
  - Einfügen eines neuen Elements mit beliebiger Priorität

## Aus der Definition folgt, dass

- ein Teilbaum eines Heaps wiederum ein Heap ist
- in einem beliebigen Pfad im Heap die Elemente sortiert sind
- Je nach dem, ob bei der Sortierung die Eigenschaft *größer* oder *kleiner* betrachtet wird, spricht man von *minHeap* oder *maxHeap*



# Priority Queues

Typische Operationen auf Priority Queues:

- isempty, insert, readmin, delmin

Für das **Einfügen**:

- hänge das einzufügende Element an den freien Knoten der vorletzten Ebene
- stelle die Heap-Eigenschaft wieder her

Für das **Entnehmen** eines minimalen/maximalen Elements:

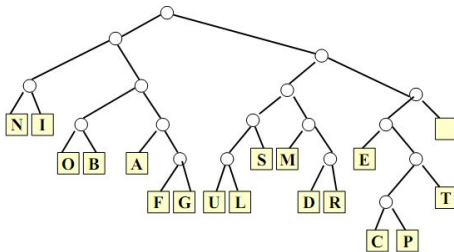
- lösche das Wurzelement
- füge das letzte Blatt an der Wurzel ein
- stelle die Heap-Eigenschaft wieder her (*heapify*)

**Aufwand:**

- Min-Lesen  $O(1)$
- Min-Löschen  $O(\log n)$
- Einfügen  $O(\log n)$

# Trie für die Huffman-Codierung von "A SIMPLE STRING ..."

	_	A	B	C	D	E	F	G	I	L	M	N	O	P	R	S	T	U
$k$	0	1	2	3	4	5	6	7	9	12	13	14	15	16	18	19	20	21
count[k]	11	3	3	1	2	5	1	2	6	2	4	5	3	1	2	4	3	2

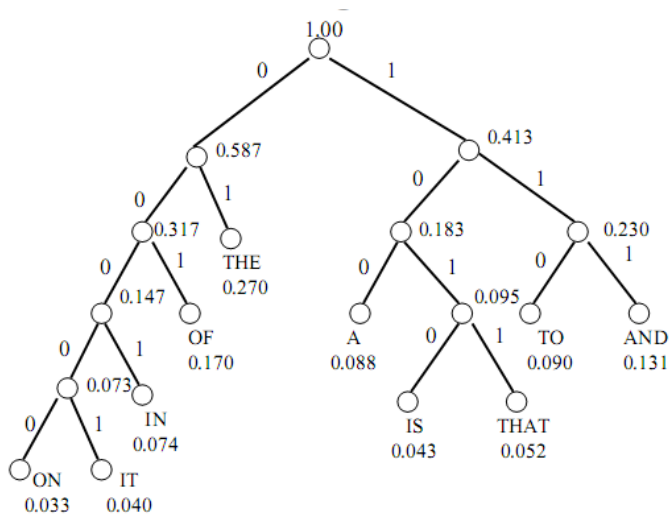


Ableitung des eigentlichen Codes aus dem Trie: Kantenlabel 0 nach links und 1 nach rechts, lese Code eines Zeichens auf Pfad von Wurzel zu Zeichen ab.

# Huffman-Codierung für Wörter der Englischen Sprache

Codierungs-Einheit	Wahrscheinlichkeit des Auftretens	Code-Wert	Code-Länge
the	0.270	01	2
of	0.170	001	3
and	0.137	111	3
to	0.099	110	3
a	0.088	100	3
in	0.074	0001	4
that	0.052	1011	4
is	0.043	1010	4
it	0.040	00001	5
on	0.033	00000	5

# Huffman-Codierungs-Baum für englische Wörter



# Häufigkeitsverteilung deutscher Bigramme

<i>en</i>	4.47	<i>nw</i>	0.55	<i>ab</i>	0.25	<i>nm</i>	0.16
<i>er</i>	3.40	<i>us</i>	0.54	<i>il</i>	0.25	<i>pe</i>	0.16
<i>ch</i>	2.80	<i>nn</i>	0.53	<i>mm</i>	0.25	<i>rl</i>	0.16
<i>nd</i>	2.58	<i>nt</i>	0.52	<i>nz</i>	0.25	<i>sm</i>	0.16
<i>ei</i>	2.26	<i>ta</i>	0.51	<i>sg</i>	0.25	<i>sp</i>	0.16
<i>de</i>	2.14	<i>eg</i>	0.50	<i>sw</i>	0.25	<i>th</i>	0.16
<i>in</i>	2.04	<i>eh</i>	0.50	<i>rn</i>	0.24	<i>wo</i>	0.16
<i>es</i>	1.81	<i>zu</i>	0.50	<i>ro</i>	0.24	<i>af</i>	0.15
<i>te</i>	1.78	<i>al</i>	0.49	<i>ea</i>	0.23	<i>lu</i>	0.15
<i>ie</i>	1.76	<i>ed</i>	0.48	<i>fr</i>	0.23	<i>mu</i>	0.15
<i>un</i>	1.73	<i>ru</i>	0.48	<i>sd</i>	0.23	<i>no</i>	0.15
<i>ge</i>	1.68	<i>rs</i>	0.47	<i>tt</i>	0.23	<i>nv</i>	0.15
<i>st</i>	1.24	<i>ig</i>	0.45	<i>tw</i>	0.23	<i>rf</i>	0.15
<i>ic</i>	1.19	<i>ts</i>	0.45	<i>gr</i>	0.22	<i>ut</i>	0.15
<i>he</i>	1.17	<i>ma</i>	0.43	<i>tz</i>	0.22	<i>br</i>	0.14
<i>ne</i>	1.17	<i>sa</i>	0.43	<i>fe</i>	0.21	<i>ez</i>	0.14
<i>se</i>	1.17	<i>wa</i>	0.43	<i>gt</i>	0.21	<i>ho</i>	0.14
<i>ng</i>	1.07	<i>ac</i>	0.42	<i>rh</i>	0.21	<i>ka</i>	0.14
<i>re</i>	1.07	<i>eu</i>	0.42	<i>ds</i>	0.20	<i>os</i>	0.14
<i>au</i>	1.04	<i>so</i>	0.41	<i>du</i>	0.20	<i>bl</i>	0.13
<i>di</i>	1.02	<i>ar</i>	0.40	<i>mi</i>	0.20	<i>dw</i>	0.13
<i>be</i>	0.96	<i>tu</i>	0.40	<i>nb</i>	0.20	<i>ep</i>	0.13
<i>ss</i>	0.94	<i>ck</i>	0.37	<i>nk</i>	0.20	<i>hm</i>	0.13

# Häufigkeitsverteilung deutscher Trigramme

<i>ein</i>	1.22	<i>ese</i>	0.27	<i>hre</i>	0.18	<i>nne</i>	0.14	<i>auc</i>	0.11
<i>ich</i>	1.11	<i>auf</i>	0.26	<i>hei</i>	0.18	<i>nes</i>	0.14	<i>als</i>	0.11
<i>nde</i>	0.89	<i>ben</i>	0.26	<i>lei</i>	0.18	<i>ond</i>	0.14	<i>alt</i>	0.11
<i>die</i>	0.87	<i>ber</i>	0.26	<i>nei</i>	0.18	<i>oen</i>	0.14	<i>eic</i>	0.11
<i>und</i>	0.87	<i>eit</i>	0.26	<i>nau</i>	0.18	<i>sdi</i>	0.14	<i>esc</i>	0.11
<i>der</i>	0.86	<i>ent</i>	0.26	<i>sge</i>	0.18	<i>sun</i>	0.14	<i>enh</i>	0.11
<i>che</i>	0.75	<i>est</i>	0.26	<i>tte</i>	0.18	<i>von</i>	0.14	<i>eil</i>	0.11
<i>end</i>	0.75	<i>sei</i>	0.26	<i>wei</i>	0.18	<i>bei</i>	0.13	<i>fen</i>	0.11
<i>gen</i>	0.71	<i>and</i>	0.25	<i>abe</i>	0.17	<i>chl</i>	0.13	<i>gan</i>	0.11
<i>sch</i>	0.66	<i>ess</i>	0.25	<i>chd</i>	0.17	<i>chn</i>	0.13	<i>hte</i>	0.11
<i>cht</i>	0.61	<i>ann</i>	0.24	<i>des</i>	0.17	<i>chw</i>	0.13	<i>iea</i>	0.11
<i>den</i>	0.57	<i>esi</i>	0.24	<i>nte</i>	0.17	<i>ech</i>	0.13	<i>ieb</i>	0.11
<i>ine</i>	0.53	<i>ges</i>	0.24	<i>rge</i>	0.17	<i>edi</i>	0.13	<i>nli</i>	0.11
<i>nge</i>	0.52	<i>nsc</i>	0.24	<i>tes</i>	0.17	<i>enk</i>	0.13	<i>rda</i>	0.11
<i>nun</i>	0.48	<i>nwi</i>	0.24	<i>uns</i>	0.17	<i>eun</i>	0.13	<i>rsc</i>	0.11
<i>ung</i>	0.48	<i>tei</i>	0.24	<i>vor</i>	0.17	<i>enz</i>	0.13	<i>std</i>	0.11
<i>das</i>	0.47	<i>eni</i>	0.23	<i>dem</i>	0.16	<i>hau</i>	0.13	<i>sst</i>	0.11
<i>hen</i>	0.47	<i>ige</i>	0.23	<i>hin</i>	0.16	<i>ite</i>	0.13	<i>tre</i>	0.11
<i>ind</i>	0.46	<i>aen</i>	0.22	<i>her</i>	0.16	<i>ief</i>	0.13	<i>uss</i>	0.11
<i>enw</i>	0.45	<i>era</i>	0.22	<i>lle</i>	0.16	<i>imm</i>	0.13	<i>all</i>	0.10
<i>ens</i>	0.44	<i>ern</i>	0.22	<i>nan</i>	0.16	<i>ihr</i>	0.13	<i>aft</i>	0.10
<i>ies</i>	0.44	<i>rde</i>	0.22	<i>tda</i>	0.16	<i>iss</i>	0.13	<i>bes</i>	0.10
<i>ste</i>	0.44	<i>ren</i>	0.22	<i>tel</i>	0.16	<i>kei</i>	0.13	<i>dei</i>	0.10
<i>ten</i>	0.44	<i>tun</i>	0.22	<i>ueb</i>	0.16	<i>mei</i>	0.13	<i>erf</i>	0.10
<i>ere</i>	0.43	<i>ing</i>	0.21	<i>ang</i>	0.15	<i>nsi</i>	0.13	<i>ess</i>	0.10
<i>lic</i>	0.42	<i>sta</i>	0.21	<i>cha</i>	0.15	<i>nem</i>	0.13	<i>esw</i>	0.10



# Künstliche Sprache (nach Kumpfmüller)

## *Einergruppen (Buchstabhäufigkeit)*

EME GKNEET ERS TITBL BTZENFNDGBGD EAI E LASZ  
BETEATR IASMIRCH EGEOM

## *Zweiergruppen (Paarhäufigkeit)*

AUSZ KEINU WONDINGLIN DUFRN ISAR STEISBERER ITEHM  
ANORER

## *Dreiergruppen*

PLANZEUNDGES PHIN INE UNDEN ÜBBEICHT GES AUF ES SO  
UNG GAN DICH WANDERSO

## *Vierergruppen*

ICH FOLGEMÄSZIG BIS STEHEN DISPONIN SEELE NAMEN

- Welche Arten von Redundanz können Lauflängencodierung und Huffmancodierung ausnutzen bzw. nicht ausnutzen?
- Welches Verfahren sollte wann benutzt werden?
- Welches Verfahren ist gut für Text?
- Sind Kombinationen, d.h. die Hintereinanderausführung von Kompressions-Verfahren sinnvoll?
- Verschiedene Kompressionsverfahren arbeiten unterschiedlich gut bei Daten unterschiedlichen Typs (Bilder, numerische Daten, Text).
- Schlussfolgerung: Wir interessieren uns für weitere Verfahren.

# Die LZ-Familie von Kompressionsalgorithmen

## LZ77, LZ78, LZW, ...

Nutze unterschiedliche Häufigkeit von *Zeichenfolgen* aus.

**Grundlegende Idee:** verwende ein “Wörterbuch” von Zeichenfolgen, so dass häufige Zeichenfolgen durch Verweise auf das Wörterbuch platzsparend codiert werden können. (LZ78,LZW: explizites Wörterbuch; LZ77: implizit, “sliding window”)

**Text-Beispiel:** RABARBABARBARABARBARENBART<sup>1</sup>

**Grundsätzlicher Konflikt:** Je grösser das Wörterbuch, desto

- + mehr und längere Zeichenfolgen können durch Verweise codiert werden
- mehr Platzbedarf pro Verweis
- mehr Platzbedarf für Wörterbuch
- höhere Laufzeit und Platzbedarf für Kompression und Dekompression

---

<sup>1</sup>in kompressionsfreundlich-reformierter Rechtschreibung

# Historischer Exkurs (LZ-Familie)

1977	Abraham Lempel und Jacob Ziv erfinden den Kompressionsalgorithmus <a href="#">LZ77</a>
1983	Terry A. Welch (Sperry Corporation - später Unisys) patentiert <a href="#">LZW</a> (Variante von LZ78)
1987	CompuServe veröffentlicht <a href="#">GIF</a> als freie und offene Spezifikation (Version <a href="#">87a</a> )
1989	Vorstellung von <a href="#">GIF 89a</a>
1993	Unisys informiert CompuServe über die Verwendung ihres patentierten LZW-Algorithmus in GIF
1994	Unisys gibt öffentlich bekannt, Gebühren für die Verwendung des LZW-Algorithmus einzufordern
1995	die <a href="#">PNG</a> Gruppe wird gegründet, erste PNG Bilder werden ins Netz gestellt, die PNG-Spezifikation 0.92 steht im <a href="#">W3C</a>
1997	PNG-Unterstützung in Netscape 4.04 und IE 4.0