

# Zeichendarstellung mit Unicode und UTF-8

## Unicode

Um Probleme, die sich zum einen mit unterschiedlichen Zeichensätzen, zum anderen auch durch andere Sprachen, die mehr als 128 Zeichen umfassen, ergeben haben, wurde der Unicode-Standard entwickelt.

Unicode ist ein internationaler Standard, in dem langfristig für jedes sinnvolle Schriftzeichen oder Textelement aller bekannten Schriftkulturen und Zeichensysteme ein digitaler Code festgelegt wird. Ziel ist es, die Verwendung unterschiedlicher und inkompatibler Kodierungen in verschiedenen Ländern oder Kulturkreisen zu beseitigen. Unicode wird ständig um Zeichen weiterer Schriftsysteme durch das Unicode-Konsortium ergänzt. (Wikipedia, <https://de.wikipedia.org/wiki/Unicode>)

Im Unicode Standard hat jedes Zeichen einen eigenen "Unicode-Code", damit lassen sich derzeit 1.111.998 elementare Zeichen ("Codepunkte") abbilden. Darstellung: U+00DF (Mindestens 4x4Bit, bis zu U+10FFFF) Diese Codepunkte bilden den Unicode Zeichensatz.

Der Codepoint eines Unicode-Zeichens ist nur eine abstrakte Nummer. Die Schreibweise dieser Nummer im Unicode-Standard erfolgt in hexadezimalen Zahlen mit vorangestelltem "U+". Der Codepoint legt noch keine computerkompatible Darstellung fest, dies ist Aufgabe der eigentlichen Codierung. Da die Unicode-Codepoints von U+0000 bis U+10FFFF (hexadezimale Zahlendarstellung), mit einer beabsichtigten Lücke zwischen U+D7FF und U+E000, reichen, sind für eine vollständige Codierung des gesamten Codepoint-Bereichs als Binärzahl mindestens 3 Byte erforderlich.<sup>1)</sup>

Frage: Wie kann man die "Zeichennummern", also die Codepunkte als Bitmuster darstellen? Dieser Vorgang ist die **Codierung** der Unicode Zeichen.

UTF-32	UTF-16	UTF-8
Jeder Codepoint wird durch 4 Bytes dargestellt. Das reicht für alle Codepoints verschwendet aber Ressourcen - einfach aber verschwenderisch.	Jeder Codepoint wird durch 2 Bytes dargestellt, halbiert also den Speicherverbrauch von UTF-32. Allerdings reichen 2 Bytes nicht für alle Codepoints, deswegen muss man tricksen, indem man die oben erwähnte Lücke zwischen U+D7FF und U+E000 nutzt. <a href="#">Details finden sich hier.</a>	UTF-8 nutzt Codiereinheiten, die 8 Bit = 1 Byte groß sind. Wie bei UTF-16 ist der Gesamtbereich der Unicode-Zeichen mit nur einem Byte nicht abzudecken, also werden bei Bedarf mehrere Bytes zusammengefasst, um einen Codepoint zu codieren. UTF-8 verwendet also je nach Bedarf zwischen 1 und 4 Byte, um einen Codepoint darzustellen.

Die Zeichen des Zeichensatzes werden wiederum auf unterschiedliche Weisen codiert, beispielsweise in Betriebssystemen. Wir betrachten die UTF-8 Kodierung von Unicode Zeichen genauer.

## UTF-8 Implementation des Unicode Zeichensatzes

Hier kann ein einzelnes Zeichen in der UTF-8-Codierung bis zu 4 Bytes umfassen, nach folgenden Regeln:

- Ist die Binärdarstellung des Unicode-Codes nicht länger als ein Byte und das erste Bit eine 0, werden die restlichen 7 Bit gemäß des ASCII Codes verwendet, die 128 verbleibenden Möglichkeiten entsprechen also genau dem ASCII-Code.
- Ist die Binärdarstellung des Unicode-Codes länger als ein Byte oder der Code ist ein Byte lang und beginnt mit einer 1 geht man wie folgt vor: Der Unicode-Code wird in 6 Bit lange Teile aufgeteilt. Für jedes dieser 6 Bit Pakete wird ein Byte zur Darstellung verwendet, jedes Byte beginnt mit '10'. Das erste Byte beginnt mit einer '1' für jedes Byte, das verwendet wird. Benötigt man also 3 Byte, um ein Zeichen in UTF-8 darzustellen, beginnt das erste Byte mit '111'. Bevor die Nutzdaten beginnen, muss noch eine Null eingefügt werden <sup>2)</sup>

### Beispiele:

(1)

$$y = 79_{16} = 0111\ 100_2$$

Beginnt mit einer Null und ist nicht länger als ein Byte → die letzten 7 Bit werden verwendet, um zu codieren, also ein "ASCII k" in UTF-8

**UTF-8: 0110 1011**

(2)

$$\text{ä} = E4_{16} = 1110\ 0100_2$$

Nur ein Byte lang, beginnt aber mit einer 1. Die 8 Bit müssen in 6 Bit Abschnitte geteilt und auf 2 Byte verteilt werden, beginnend von rechts, links wird stets mit 0en aufgefüllt:

000011 100100

- Das zweite Byte beginnt nach den Regeln mit 10, daran schließen die Nutzdaten an: 10 100100
- Das erste Byte beginnt mit 11 (weil man zwei Byte benötigt) dann wird mit 0en aufgefüllt, dann kommen die Nutzdaten: 11 000011

Die UTF-8 Codierung des Unicode-ä ist also 1100 0011 1010 0100. Die Nutzdaten, die den Code des Unicode Zeichens transportieren sind in jeden Byte nur die jeweils letzten 6 Bit.

(3)

# 乐

Chinese, Japanese, Korean (CJK) unified ideograph (U+4E50)

乐  $\rightarrow$  U+4E50  $\rightarrow$  4E50<sub>16</sub>  $\rightarrow$  0100 1110 0101 0000<sub>2</sub>

- 16 Bit Daten zu codieren, dafür braucht man 3 Byte (  $3 \times 6 = 18$  )
- Der UTF-8 Code beginnt also mit der Startsequenz 1110
- Dann von rechts beginnend 6 Bit (01 000), das Byte beginnt mit 10 (Regel) also ist das dritte Byte 1001 1000
- Die nächsten 6 Bit analog: 1110 01  $\rightarrow$  1011 1001
- Die fehlenden 4 Bit 0100 mit Padding + Startsequenz (111) ergeben das erste Byte 1110 0100

Die UTF-8 Codierung des Unicode-Zeichens 乔 ist also 3 Byte lang und sieht so aus: 1110 0100  
1011 1001 1001 0000



## (A3)

Wandle die nachfolgenden Zeichen des Unicode Zeichensatzes in die UTF-8-Codierung um. Der Hexadezimalcode des Unicode Zeichens ist jeweils angegeben.

Gehe jeweils wie in den Beispielen oben vor. Markiere die "Nutzdaten" die das eigentlich Unicode-Zeichen "transportieren".

1. I=49<sub>16</sub>
2. Ö=D6<sub>16</sub>
3. 弈=5F08<sub>16</sub>
4. □=1F60A<sub>16</sub>

## Lösung 1

0**1001001**, ein Byte, erstes Bit 0.

## Lösung 2

110**00011** 100**10110**

## Lösung 3

11100101 10111100 10001000

## Lösung 4

11110000 10011111 10011000 10001010

---



### (A4)

Wie viele unterschiedliche Unicode-Zeichen lassen sich theoretisch mit 1 Byte, 2 Bytes, 3 Bytes und 4 Bytes unter Beachtung der UTF-8-Regeln darstellen?

### Lösung

- 1 Byte: 7 nutzbare Bits  $\rightarrow 2^7 = 128$  Zeichen
- 2 Bytes: 5+6 = 11 nutzbare Bits  $\rightarrow 2^{11} = 2\,048$  Zeichen
- 3 Bytes: 4+6+6 = 16 nutzbare Bits  $\rightarrow 2^{16} = 65\,536$  Zeichen
- 4 Bytes: 3+6+6+6 = 21 nutzbare Bits  $\rightarrow 2^{21} = 2\,097\,152$  Zeichen

## Material

<a href="#">03-zeichencodierung.odp</a>	3.7 MiB 22.09.2022 07:31
<a href="#">03-zeichencodierung.pdf</a>	1.9 MiB 22.09.2022 07:31
<a href="#">asciitable.png</a>	724.3 KiB 13.10.2021 13:48
<a href="#">cs.png</a>	25.6 KiB 15.10.2021 10:10

---

CC-BY-SA Frank Schiebel, mit Material von Kimmig, ZPG Informatik BW

<sup>1)</sup> Einzelne Teile dieser Wikiseite sind dem SelfHTML Wiki entnommen und stehen unter CC-BY-SA Lizenz <https://wiki.selfhtml.org/wiki/Zeichencodierung/Unicode>

<sup>2)</sup> Warum?

From:  
<https://www.info-bw.de/> -

Permanent link:  
<https://www.info-bw.de/faecher:informatik:oberstufe:codierung:utf8:start?rev=1701166302>

Last update: **28.11.2023 10:11**

