

Buchstaben, Zeichen und Text

Informatik Grundlagen
Kantonsschule am Burggraben

Ivo Blöchliger

Alles Bits und Bytes

- Der Computer speichert nur Bits und Bytes
- Im Prinzip alles als natürliche Zahl interpretierbar.
 - z.B. ist ein Video eine einzige, riesige Zahl.
- Codierung:
 - Abmachung, wie die Bytes zu interpretieren sind.

Buchstaben

- Idee: Buchstaben (und andere Zeichen) nummerieren.
- Standard ASCII (7 Bits!)
- Was ergeben folgende Bytes, wenn man diese als ASCII-Zeichen interpretiert?

0x4f 0x6b 0x21 0x0a 0x36 0x2e 0x30

- Suchen Sie dazu eine ASCII-Tabelle im Internet

0x4f 0x6b 0x21 0x0a 0x36 0x2e 0x30

0x4f 0x6b 0x21 0x0a 0x36 0x2e 0x30

ASCII-Zeichentabelle, hexadezimale Nummerierung

	...0	...1	...2	...3	...4	...5	...6	...7	...8	...9	...A	...B	...C	...D	...E	...F
0...	<i>NUL</i>	<i>SOH</i>	<i>STX</i>	<i>ETX</i>	<i>EOT</i>	<i>ENQ</i>	<i>ACK</i>	<i>BEL</i>	<i>BS</i>	<i>HT</i>	<i>LF</i>	<i>VT</i>	<i>FF</i>	<i>CR</i>	<i>SO</i>	<i>SI</i>
1...	<i>DLE</i>	<i>DC1</i>	<i>DC2</i>	<i>DC3</i>	<i>DC4</i>	<i>NAK</i>	<i>SYN</i>	<i>ETB</i>	<i>CAN</i>	<i>EM</i>	<i>SUB</i>	<i>ESC</i>	<i>FS</i>	<i>GS</i>	<i>RS</i>	<i>US</i>
2...	<i>SP</i>	<i>!</i>	<i>"</i>	<i>#</i>	<i>\$</i>	<i>%</i>	<i>&</i>	<i>'</i>	<i>(</i>	<i>)</i>	<i>*</i>	<i>+</i>	<i>,</i>	<i>-</i>	<i>.</i>	<i>/</i>
3...	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>:</i>	<i>;</i>	<i><</i>	<i>=</i>	<i>></i>	<i>?</i>
4...	<i>@</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>	<i>G</i>	<i>H</i>	<i>I</i>	<i>J</i>	<i>K</i>	<i>L</i>	<i>M</i>	<i>N</i>	<i>O</i>
5...	<i>P</i>	<i>Q</i>	<i>R</i>	<i>S</i>	<i>T</i>	<i>U</i>	<i>V</i>	<i>W</i>	<i>X</i>	<i>Y</i>	<i>Z</i>	<i>[</i>	<i>\</i>	<i>]</i>	<i>^</i>	<i>_</i>
6...	<i>`</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>	<i>i</i>	<i>j</i>	<i>k</i>	<i>l</i>	<i>m</i>	<i>n</i>	<i>o</i>
7...	<i>p</i>	<i>q</i>	<i>r</i>	<i>s</i>	<i>t</i>	<i>u</i>	<i>v</i>	<i>w</i>	<i>x</i>	<i>y</i>	<i>z</i>	<i>{</i>	<i> </i>	<i>}</i>	<i>~</i>	<i>DEL</i>

Was ist mit ä, ö, ù?

- ASCII-Code nur lateinische Buchstaben ohne Akzente
- Nur 7 Bits, also nochmals 128 Zeichen mit achtem Bit!
- Früher (und Windows noch heute):
 - Viele mögliche Zeichentabelle, je nach Sprachregion
 - Bei uns gebräuchlich: Latin-1 (ISO8859-1), CP-1252
- Heute (bis auf Windows fast ausschliesslich):

Unicode: Jedem Zeichen eine Nummer

Momentan über 100'000 Zeichen spezifiziert.

Unicode

- Weltweit jedem Zeichen seine Nummer
- Emojis, Akzente, Schriftrichtung, Ligaturen
 - Siehe Unicode-Video
- Codierung meist UTF-8:
 - Standard-ASCII mit 1 Byte
 - Andere Zeichen mit 2 oder mehr Bytes
 - Erstes Byte gibt auch an, aus wie vielen Bytes das Zeichen besteht.
 - Speichereffizient für Sprachen mit lateinischem Alphabet.

Text-Dateien (nicht Word!)

- Folge von Buchstaben (Bytes)
 - Und Steuerzeichen wie Zeilenumbrüche ‘\n’, Tabulatoren ‘\t’,...
- Universal lesbar und veränderbar, wenn in ASCII
 - Unicode in UTF-8 codiert heute ebenfalls.
- Keine Stil-Information (Schriftart, Grösse, Farbe, etc.)
- Beispiel: Python-Code, CSV-Dateien (Tabellarische Daten)

Und mit Stil?

- Variante 1
Proprietäres Format, nur mit einem Programm lesbar
- Variante 2
Markup: z.B. HTML, XML
`<h1>Titel</h1>Rot!`
- Office-Dokumente heute: zip-Archiv von XML-Dateien

Und in Python?

Und in Python?

```
>>> ord("A")
```

```
65
```

```
>>> chr(97)
```

```
'a'
```

```
>>> s = chr(0x1f600)
```

```
>>> s
```

```
⋮
```

```
>>> [bin(x) for x in s.encode("utf-8")]
```

```
['0b11110000', '0b10011111', '0b10011000', '0b10000000']
```