

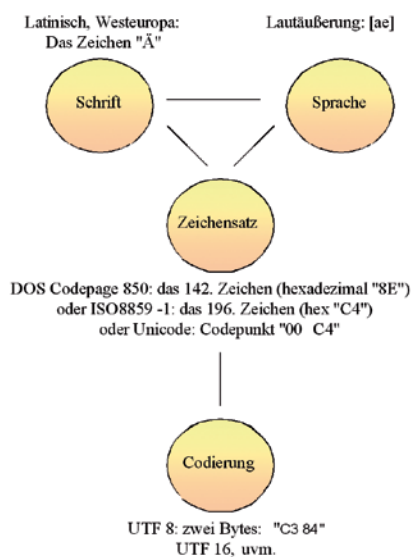
also 255 Positionen, zur Verfügung stehen und deshalb der Druck zur Platzeinsparung gering ist.

Im Unicode aber gibt es eine Hintertür für Nutzer des ASCII-Zeichensatzes, also der lateinischen Schrift ohne Umlaute und Sonderzeichen. Diese Codepunkte sind so angelegt, dass man sie trotzdem mit nur einem Byte darstellen kann. Sie liegen nämlich zwischen 0000 und 007F, also zwischen 0 und 127, wie bisher auch. Auf diese Weise ist ASCII kompatibel zu Unicode.

Erst für andere Schriften muss man dann weitere Bits und zusätzlich auch Zählmarker, Startbits und Stopbits hinzunehmen, der Codepunkt benötigt dann zwei oder sogar drei oder vier Bytes. Diese Technik der Formatierung nennt man UTF-8 (8-bit unicode transformation format). Und sie ist auch am weitesten verbreitet.

Eine weitere, und ältere Möglichkeit ist die UTF-16 Codierung: Hier werden grundsätzlich immer zwei Bytes benutzt, so dass die Codepoints von 0000 bis FFFF vollständig abgedeckt sind. Für darüberliegende Codepoints muss man dann allerdings zwei „2-Byte“-Zeichen kombinieren.

Die Codierung kann auch als Formatierung oder „Encoding“ bezeichnet werden. Der Zusammenhang zwischen Schrift, Zeichensatz und Codierung ist in folgendem Schaubild dargestellt, und zwar für das Zeichen „Ä“, welches nicht zu den ersten 127 Codepunkten gehört.



Der Zusammenhang zwischen Sprache, Schrift, Zeichensatz und Codierung

Probleme

Damit sind aber bereits die ersten Nachteile des Unicode offenbar: Da man normalerweise mehr als ein Byte zur Codierung benötigt, nimmt die Menge an Speicherplatz zu, die man für einen Text verbraucht. Zum Beispiel kommt man bei einem herkömmlichen Zeichensatz wie ISO 8859-5 mit einem Byte für kyrillische Buchstaben aus, weil man ja alle anderen Schriften in solchen Zeichensätzen ausblendet. Im Unicode erreicht man hingegen leicht eine Verdoppelung der Datenmenge. Für größere Datenbanken wie zum Beispiel landesweite Straßennetze mit Namen kann das schnell eine kritische Schwelle überschreiten.

Infolgedessen kann auch die Leistungsfähigkeit von Anwendungen leiden. Es muss ja unter Umständen die doppelte Datenmenge gelesen werden, um dieselben Zeichen auf dem Bildschirm darzustellen. Man sollte also gut überlegen, ob für die benötigten Daten der Einsatz von Unicode überhaupt gerechtfertigt ist.

Schließlich muss das Betriebssystem auch eine Schriftart bereitstellen, die den Unicode für uns lesbar machen kann. Wer mit einem neueren Windows-Rechner arbeitet, findet in den Anwendungen in der Regel Schriftarten wie „Lucida Sans Unicode“ oder „Arial Unicode“, die einen großen Teil des Problems aus der Welt schaffen und für normale Anwendungsfälle recht gut brauchbar sind.

Auf anderen Betriebssystemen oder älteren Rechnern muss das aber nicht der Fall sein. Und wer auf weniger triviale Schriftarten zurückgreifen möchte, hat dazu in der Regel keine Unicode-fähige Variante und muss dann ohnehin einen passenden anderen Font kaufen.

Internationalisierung

Dennoch – die Zeichen der Zeit stehen auf Internationalisierung. Vorbei ist die Vorherrschaft der herkömmlichen IT-Welt, die mit 26

Buchstaben auskommt. Längst mischen die Menschen aus Asien und Osteuropa kräftig mit in den Märkten, bauen und nutzen Software und Internet. Und sie möchten natürlich nicht auf die Besonderheiten ihrer Schrift verzichten. Wer sich einmal über „München“ oder „Munchn“ auf ausländischen Internetseiten amüsiert hat, der weiß, wie sich Osteuropäer fühlen, die mit Schreibfehlern in ihren Ortsnamen konfrontiert werden. Mit Unicode hingegen können verschiedene Schriften gleichzeitig dargestellt werden.

Und damit wird das Beherrschen der Schriften und Symbole dieser Welt zum Erfolgsfaktor auf dem Markt. Denn wer „Łodz“ in seinen Daten richtig schreibt, dem traut der Kunde mehr zu als dem Urheber von „Łodz“. Die identitätsstiftende Sprach- und Schriftkultur ist ein so grundlegender Teil von uns selbst, dass sie oft unbemerkt kaufentscheidend sein kann. Der internationale Datenaustausch ist heute ohne die Vorteile von Unicode kaum noch denkbar.



Mehrere Schriftarten in derselben Karte – mit Unicode kein Problem

In Kürze

Um Zeichen zu sparen, kann man es in den Worten des Unicode-Konsortiums kurz auf den Punkt bringen: „Unicode bietet eine eindeutige Nummer für jedes Zeichen, egal welche Plattform, egal welches Programm, egal welche Sprache.“ (<http://www.unicode.org/>).



Dipl.-Geoökol. Thomas Engel, Jahrgang 1975, studierte an den Universitäten Karlsruhe und Hamburg Geoökologie mit Schwerpunkt Geoinformatik. Seit acht Jahren arbeitet er bei der PTV AG und betreut zahlreiche Datenversorgungsprojekte für alle geographischen Applikationen des Konzerns. Schwerpunkte sind das Standardformat GDF, die Umsetzung internationaler Geodaten und die Konzeptionierung der Mautberechnung.