

The Unicode HOWTO

Table of Contents

The Unicode HOWTO	1
<u>Bruno Haible, <haible@clisp.cons.org></u>	1
<u>1. Introduction</u>	1
<u>1.1 Why Unicode?</u>	1
<u>1.2 Unicode encodings</u>	1
<u>Footnotes for C/C++ developers</u>	2
<u>1.3 Related resources</u>	3
<u>2. Display setup</u>	3
<u>2.1 Linux console</u>	3
<u>2.2 X11 Foreign fonts</u>	4
<u>2.3 X11 Unicode fonts</u>	5
<u>2.4 Unicode xterm</u>	5
<u>2.5 TrueType fonts</u>	6
<u>2.6 Miscellaneous</u>	8
<u>3. Locale setup</u>	8
<u>3.1 Files & the kernel</u>	8
<u>3.2 Upgrading the C library</u>	8
<u>3.3 General data conversion</u>	9
<u>3.4 Locale environment variables</u>	9
<u>3.5 Creating the locale support files</u>	10
<u>4. Specific applications</u>	10
<u>4.1 Shells</u>	10
<u>bash</u>	10
<u>4.2 Networking</u>	11
<u>telnet</u>	11
<u>kermi</u>	11
<u>4.3 Browsers</u>	11
<u>Netscape</u>	11
<u>Mozilla</u>	11
<u>Amaya</u>	12
<u>lynx</u>	12
<u>w3m</u>	13
<u>Test pages</u>	13
<u>4.4 Editors</u>	13
<u>yudit</u>	13
<u>yudit-1.5</u>	13
<u>yudit-2.1</u>	14
<u>Fonts for yudit</u>	14
<u>vim</u>	14
<u>cooledit</u>	14
<u>emacs</u>	15
<u>xemacs</u>	17
<u>nedit</u>	18
<u>xedit</u>	18
<u>axe</u>	18
<u>pico</u>	18
<u>mined98</u>	18
<u>qemacs</u>	18

Table of Contents

The Unicode HOWTO

<u>4.5 Mailers</u>	19
<u>pine</u>	19
<u>kmail</u>	19
<u>Netscape Communicator</u>	19
<u>emacs (rmail, vm)</u>	20
<u>mutt</u>	20
<u>exmh</u>	20
<u>4.6 Text processing</u>	20
<u>groff</u>	20
<u>TeX</u>	20
<u>4.7 Databases</u>	21
<u>PostgreSQL</u>	21
<u>Interbase</u>	21
<u>4.8 Other text-mode applications</u>	21
<u>less</u>	21
<u>lv</u>	21
<u>expand</u>	21
<u>col, colcrt, colrm, column, rev, ul</u>	21
<u>figlet</u>	22
<u>Base utilities</u>	22
<u>4.9 Other X11 applications</u>	29
<u>5. Printing</u>	29
<u>5.1 Printing using TrueType fonts</u>	29
<u>uniprint</u>	29
<u>wprint</u>	29
<u>Comparison</u>	30
<u>5.2 Printing using fixed-size fonts</u>	30
<u>txtbdf2ps</u>	30
<u>5.3 The classical approach</u>	30
<u>TeX, Omega</u>	30
<u>DocBook</u>	30
<u>groff -Tps</u>	30
<u>5.4 No luck with</u>	30
<u>Netscape's "Print..."</u>	31
<u>Mozilla's "Print..."</u>	31
<u>html2ps</u>	31
<u>a2ps</u>	31
<u>enscript</u>	31
<u>6. Making your programs Unicode aware</u>	31
<u>6.1 C/C++</u>	31
<u>For normal text handling</u>	31
<u>Portability notes</u>	32
<u>The libutf8 library</u>	34
<u>The Plan9 way</u>	34
<u>For graphical user interface</u>	34
<u>For advanced text handling</u>	34
<u>For conversion</u>	35

Table of Contents

The Unicode HOWTO

<u>iconv</u>	35
<u>librecode</u>	35
<u>ICU</u>	36
<u>Other approaches</u>	36
<u>6.2 Java</u>	36
<u>6.3 Lisp</u>	37
<u>6.4 Ada95</u>	37
<u>6.5 Python</u>	37
<u>6.6 JavaScript/ECMAScript</u>	38
<u>6.7 Tcl</u>	38
<u>6.8 Perl</u>	38
<u>6.9 Related reading</u>	38
<u>7. Other sources of information</u>	38
<u>7.1 Mailing lists</u>	38
<u>linux-utf8</u>	38
<u>li18nux</u>	39
<u>unicode</u>	39
<u>X11 internationalization</u>	39
<u>X11 fonts</u>	39

The Unicode HOWTO

Bruno Haible, <haible@clisp.cons.org>

v1.0, 23 January 2001

This document describes how to change your Linux system so it uses UTF-8 as text encoding. - This is work in progress. Any tips, patches, pointers, URLs are very welcome.

1. Introduction

1.1 Why Unicode?

People in different countries use different characters to represent the words of their native languages. Nowadays most applications, including email systems and web browsers, are 8-bit clean, i.e. they can operate on and display text correctly provided that it is represented in an 8-bit character set, like ISO-8859-1.

There are far more than 256 characters in the world - think of cyrillic, hebrew, arabic, chinese, japanese, korean and thai -, and new characters are being invented now and then. The problems that come up for users are:

- It is impossible to store text with characters from different character sets in the same document. For example, I can cite russian papers in a German or French publication if I use TeX, xdvi and PostScript, but I cannot do it in plain text.
- As long as every document has its own character set, and recognition of the character set is not automatic, manual user intervention is inevitable. For example, in order to view the homepage of the XTeamLinux distribution <http://www.xteamlinux.com.cn/> I had to tell Netscape that the web page is coded in GB2312.
- New symbols like the Euro are being invented. ISO has issued a new standard ISO-8859-15, which is mostly like ISO-8859-1 except that it removes some rarely used characters (the old currency sign) and replaced it with the Euro sign. If users adopt this standard, they have documents in different character sets on their disk, and they start having to think about it daily. But computers should make things simpler, not more complicated.

The solution of this problem is the adoption of a world-wide usable character set. This character set is Unicode <http://www.unicode.org/>. For more info about Unicode, do ``man 7 unicode'` (manpage contained in the man-pages-1.20 package).

1.2 Unicode encodings

This reduces the user's problem of dealing with character sets to a technical problem: How to transport Unicode characters using the 8-bit bytes? 8-bit units are the smallest addressing units of most computers and also the unit used by TCP/IP network connections. The use of 1 byte to represent 1 character is, however, an accident of history, caused by the fact that computer development started in Europe and the U.S. where 96 characters were found to be sufficient for a long time.

There are basically four ways to encode Unicode characters in bytes:

The Unicode HOWTO

UTF-8

128 characters are encoded using 1 byte (the ASCII characters). 1920 characters are encoded using 2 bytes (Roman, Greek, Cyrillic, Coptic, Armenian, Hebrew, Arabic characters). 63488 characters are encoded using 3 bytes (Chinese and Japanese among others). The other 2147418112 characters (not assigned yet) can be encoded using 4, 5 or 6 characters. For more info about UTF-8, do `man 7 utf-8` (manpage contained in the `man-pages-1.20` package).

UCS-2

Every character is represented as two bytes. This encoding can only represent the first 65536 Unicode characters.

UTF-16

This is an extension of UCS-2 which can represent 1112064 Unicode characters. The first 65536 Unicode characters are represented as two bytes, the other ones as four bytes.

UCS-4

Every character is represented as four bytes.

The space requirements for encoding a text, compared to encodings currently in use (8 bit per character for European languages, more for Chinese/Japanese/Korean), is as follows. This has an influence on disk storage space and network download speed (when no form of compression is used).

UTF-8

No change for US ASCII, just a few percent more for ISO-8859-1, 50% more for Chinese/Japanese/Korean, 100% more for Greek and Cyrillic.

UCS-2 and UTF-16

No change for Chinese/Japanese/Korean. 100% more for US ASCII and ISO-8859-1, Greek and Cyrillic.

UCS-4

100% more for Chinese/Japanese/Korean. 300% more for US ASCII and ISO-8859-1, Greek and Cyrillic.

Given the penalty for US and European documents caused by UCS-2, UTF-16, and UCS-4, it seems unlikely that these encodings have a potential for wide-scale use. The Microsoft Win32 API supports the UCS-2 encoding since 1995 (at least), yet this encoding has not been widely adopted for documents - SJIS remains prevalent in Japan.

UTF-8 on the other hand has the potential for wide-scale use, since it doesn't penalize US and European users, and since many text processing programs don't need to be changed for UTF-8 support.

In the following, we will describe how to change your Linux system so it uses UTF-8 as text encoding.

Footnotes for C/C++ developers

The Microsoft Win32 approach makes it easy for developers to produce Unicode versions of their programs: You `#define UNICODE` at the top of your program and then change many occurrences of `'char'` to `'TCHAR'`, until your program compiles without warnings. The problem with it is that you end up with two versions of your program: one which understands UCS-2 text but no 8-bit encodings, and one which understands only old 8-bit encodings.

Moreover, there is an endianness issue with UCS-2 and UCS-4. The IANA character set registry <http://www.isi.edu/in-notes/iana/assignments/character-sets> says about ISO-10646-UCS-2: "this needs to specify network byte order: the standard does not specify". Network byte order is big endian. And RFC 2152

The Unicode HOWTO

is even clearer: "ISO/IEC 10646-1:1993(E) specifies that when characters the UCS-2 form are serialized as octets, that the most significant octet appear first." Whereas Microsoft, in its C/C++ development tools, recommends to use machine-dependent endianness (i.e. little endian on ix86 processors) and either a byte-order mark at the beginning of the document, or some statistical heuristics(!).

The UTF-8 approach on the other hand keeps ``char*'` as the standard C string type. As a result, your program will handle US ASCII text, independently of any environment variables, and will handle both ISO-8859-1 and UTF-8 encoded text provided the LANG environment variable is set accordingly.

1.3 Related resources

Markus Kuhn's very up-to-date resource list:

- <http://www.cl.cam.ac.uk/~mgk25/unicode.html>
- <http://www.cl.cam.ac.uk/~mgk25/ucs-fonts.html>

Roman Czyborra's overview of Unicode, UTF-8 and UTF-8 aware programs: <http://czyborra.com/utf/#UTF-8>

Some example UTF-8 files:

- In Markus Kuhn's ucs-fonts package: [quickbrown.txt](#), [UTF-8-test.txt](#), [UTF-8-demo.txt](#).
- <http://www.columbia.edu/kermit/utf8.html>
- <ftp://ftp.cs.su.oz.au/gary/x-utf8.html>
- The file `iso10646` in the Kosta Kostis' trans-1.1.1 package
<ftp://ftp.nid.ru/pub/os/unix/misc/trans111.tar.gz>
- <ftp://ftp.dante.de/pub/tex/info/lwc/apc/utf8.html>
- <http://www.cogsci.ed.ac.uk/~richard/unicode-sample.html>

2. Display setup

We assume you have already adapted your Linux console and X11 configuration to your keyboard and locale. This is explained in the Danish/International HOWTO, and in the other national HOWTOs: Finnish, French, German, Italian, Polish, Slovenian, Spanish, Cyrillic, Hebrew, Chinese, Thai, Esperanto. But please do not follow the advice given in the Thai HOWTO, to pretend you were using ISO-8859-1 characters (U0000..U00FF) when what you are typing are actually Thai characters (U0E01..U0E5B). Doing so will only cause problems when you switch to Unicode.

2.1 Linux console

I'm not talking much about the Linux console here, because on those machines on which I don't have `xdm` running, I use it only to type my login name, my password, and `"xinit"`.

Anyway, the `kbd-0.99` package <ftp://sunsite.unc.edu/pub/Linux/system/keyboards/kbd-0.99.tar.gz> and a heavily extended version, the `console-tools-0.2.3` package <ftp://sunsite.unc.edu/pub/Linux/system/keyboards/console-tools-0.2.3.tar.gz> contains in the `kbd-0.99/src/` (or `console-tools-0.2.3/screenfonttools/`) directory two programs: ``unicode_start'` and ``unicode_stop'`. When you call ``unicode_start'`, the console's screen output is interpreted as UTF-8. Also, the keyboard is put into Unicode mode (see `"man kbd_mode"`). In this mode, Unicode characters typed as `Alt-x1 ... Alt-xn` (where `x1,...,xn` are digits on the numeric keypad) will be emitted in UTF-8. If your keyboard or, more precisely, your

The Unicode HOWTO

normal keymap has non-ASCII letter keys (like the German Umlaute) which you would like to be CapsLockable, you need to apply the kernel patch [linux-2.2.9-keyboard.diff](#) or [linux-2.3.12-keyboard.diff](#).

You will want to use display characters from different scripts on the same screen. For this, you need a Unicode console font. The <ftp://sunsite.unc.edu/pub/Linux/system/keyboards/kbd-0.99.tar.gz> and <ftp://sunsite.unc.edu/pub/Linux/system/keyboards/console-data-1999.08.29.tar.gz> packages contain a font (LatArCyrHeb-{08,14,16,19}.psf) which covers Latin, Cyrillic, Hebrew, Arabic scripts. It covers ISO 8859 parts 1,2,3,4,5,6,8,9,10 all at once. To install it, copy it to /usr/lib/kbd/consolefonts/ and execute "/usr/bin/setfont /usr/lib/kbd/consolefonts/LatArCyrHeb-14.psf".

A more flexible approach is given by Dmitry Yu. Bolkhovityanov <D.Yu.Bolkhovityanov@inp.nsk.su> in <http://www.inp.nsk.su/~bolkhov/files/fonts/univga/index.html> and <http://www.inp.nsk.su/~bolkhov/files/fonts/univga/uni-vga.tgz>. To work around the constraint that a VGA font can only cover 512 characters simultaneously, he provides a rich Unicode font (2279 characters, covering Latin, Greek, Cyrillic, Hebrew, Armenian, IPA, math symbols, arrows, and more) in the typical 8x16 size and a script which permits to extract any 512 characters as a console font.

If you want cut&paste to work with UTF-8 consoles, you need the patch [linux-2.3.12-console.diff](#) from Edmund Thomas Grimley Evans and Stanislav Voronyi.

In April 2000, Edmund Thomas Grimley Evans <edmund@rano.org> has implemented an UTF-8 console terminal emulator. It uses Unicode fonts and relies on the Linux frame buffer device.

2.2 X11 Foreign fonts

Don't hesitate to install Cyrillic, Chinese, Japanese etc. fonts. Even if they are not Unicode fonts, they will help in displaying Unicode documents: at least Netscape Communicator 4 and Java will make use of foreign fonts when available.

The following programs are useful when installing fonts:

- "mkfontdir directory" prepares a font directory for use by the X server, needs to be executed after installing fonts in a directory.
- "xset -q | sed -e '1,/^Font Path:/d' | sed -e '2,\$d' -e 's/^ //' " displays the X server's current font path.
- "xset fp+ directory" adds a directory to the X server's current font path. To add a directory permanently, add a "FontPath" line to your /etc/XF86Config file, in section "Files".
- "xset fp rehash" needs to be executed after calling mkfontdir on a directory that is already contained in the X server's current font path.
- "xfonstsel" allows you to browse the installed fonts by selecting various font properties.
- "xlsfonts -fn fontpattern" lists all fonts matching a font pattern. Also displays various font properties. In particular, "xlsfonts -ll -fn font" lists the font properties CHARSET_REGISTRY and CHARSET_ENCODING, which together determine the font's encoding.
- "xfd -fn font" displays a font page by page.

The following fonts are freely available (not a complete list):

- The ones contained in XFree86, sometimes packaged in separate packages. For example, SuSE has only normal 75dpi fonts in the base `xf86' package. The other fonts are in the packages `xfnt100', `xfntbig', `xfntcyr', `xfntsc1'.

The Unicode HOWTO

- The Emacs international fonts, <ftp://ftp.gnu.org/pub/gnu/intlfonts/intlfonts-1.2.tar.gz> As already mentioned, they are useful even if you prefer XEmacs to GNU Emacs or don't use any Emacs at all.

2.3 X11 Unicode fonts

Applications wishing to display text belonging to different scripts (like Cyrillic and Greek) at the same time, can do so by using different X fonts for the various pieces of text. This is what Netscape Communicator and Java do. However, this approach is more complicated, because instead of working with `'Font'` and `'XFontStruct'`, the programmer has to deal with `'XFontSet'`, and also because not all fonts in the font set need to have the same dimensions.

- Markus Kuhn has assembled fixed-width 75dpi fonts with Unicode encoding covering Latin, Greek, Cyrillic, Armenian, Georgian, Hebrew scripts and many symbols. They cover ISO 8859 parts 1,2,3,4,5,7,8,9,10,13,14,15,16 all at once. These fonts are required for running xterm in utf-8 mode. They are now contained in XFree86 4.0.1, therefore you need to install them manually only if you have an older XFree86 3.x version. <http://www.cl.cam.ac.uk/~mgk25/download/ucs-fonts.tar.gz>.
- Markus Kuhn has also assembled double-width fixed 75dpi fonts with Unicode encoding covering Chinese, Japanese and Korean. These fonts are contained in XFree86 4.0.1 as well. <http://www.cl.cam.ac.uk/~mgk25/download/ucs-fonts-asian.tar.gz>
- Roman Czyborra has assembled an 8x16 / 16x16 75dpi font with Unicode encoding covering a huge part of Unicode. Download `unifont.hex.gz` and `hex2bdf` from <http://czyborra.com/unifont/>. It is not fixed-width: 8 pixels wide for European characters, 16 pixels wide for Chinese characters. Installation instructions:

```
$ gunzip unifont.hex.gz
$ hex2bdf < unifont.hex > unifont.bdf
$ bdf2pcf -o unifont.pcf unifont.bdf
$ gzip -9 unifont.pcf
# cp unifont.pcf.gz /usr/X11R6/lib/X11/fonts/misc
# cd /usr/X11R6/lib/X11/fonts/misc
# mkfontdir
# xset fp rehash
```

- Primoz Peterlin has assembled an ETL family fonts covering Latin, Greek, Cyrillic, Armenian, Georgian, Hebrew scripts. <ftp://ftp.x.org/contrib/fonts/etl-unicode.tar.gz> Use the "bdf2pcf" program in order to install it.
- Mark Leisher has assembled a proportional, 17 pixel high (12 point), font, called ClearlyU, covering Latin, Greek, Cyrillic, Armenian, Georgian, Hebrew, Thai, Laotian scripts. <http://crl.nmsu.edu/~mleisher/cu.html>. Installation instructions:

```
$ bdf2pcf -o cul2.pcf cul2.bdf
$ gzip -9 cul2.pcf
# cp cul2.pcf.gz /usr/X11R6/lib/X11/fonts/misc
# cd /usr/X11R6/lib/X11/fonts/misc
# mkfontdir
# xset fp rehash
```

2.4 Unicode xterm

xterm is part of X11R6 and XFree86, but is maintained separately by Tom Dickey.

<http://www.clark.net/pub/dickey/xterm/xterm.html> Newer versions (patch level 146 and above) contain support for converting keystrokes to UTF-8 before sending them to the application running in the xterm, and for displaying Unicode characters that the application outputs as UTF-8 byte sequence. It also contains

The Unicode HOWTO

support for double-wide characters (mostly CJK ideographs) and combining characters, contributed by Robert Brady <robert@suse.co.uk>.

To get an UTF-8 xterm running, you need to:

- Fetch <http://www.clark.net/pub/dickey/xterm/xterm.tar.gz>,
- Configure it by calling `./configure --enable-wide-chars ...`, then compile and install it.
- Have a Unicode fixed-width font installed. Markus Kuhn's `ucs-fonts.tar.gz` (see above) is made for this.
- Start `xterm -u8 -fn 'misc-fixed-medium-r-semicondensed--13-120-75-75-c-60-iso10646-1'`. The option `-u8` turns on Unicode and UTF-8 handling. The font designated by the long `-fn` option is Markus Kuhn's Unicode font. Without this option, the default font called `"fixed"` would be used, an ISO-8859-1 6x13 font.
- Take a look at the sample files contained in Markus Kuhn's `ucs-fonts` package:

```
$ cd ../ucs-fonts
$ cat quickbrown.txt
$ cat utf-8-demo.txt
```

You should be seeing (among others) greek and russian characters.

- To make xterm come up with UTF-8 handling each time it is started, add the lines

```
xterm*utf8: 1
xterm*VT100*font: -misc-fixed-medium-r-semicondensed--13-120-75-75-c-60-iso10646-1
xterm*VT100*wideFont: -misc-fixed-medium-r-normal-ja-13-125-75-75-c-120-iso10646-1
xterm*VT100*boldFont: -misc-fixed-bold-r-semicondensed--13-120-75-75-c-60-iso10646-1
```

to your `$HOME/.Xdefaults` (for yourself only). For CJK text processing with double-width characters, the following settings are probably better:

```
xterm*VT100*font: -Misc-Fixed-Medium-R-Normal--18-120-100-100-C-90-ISO10646-1
xterm*VT100*wideFont: -Misc-Fixed-Medium-R-Normal-ja-18-120-100-100-C-180-ISO10646-1
```

I don't recommend changing the system-wide `/usr/X11R6/lib/X11/app-defaults/XTerm`, because then your changes will be erased next time you upgrade to a new XFree86 version.

2.5 TrueType fonts

The fonts mentioned above are fixed size and not scalable. For some applications, especially printing, high resolution fonts are necessary, though. The most important type of scalable, high resolution fonts are TrueType fonts. They are currently supported by

- XFree86 4.0.1; you need to add the line

```
Load "freetype"
```

or

```
Load "xft"
```

to the `"Module"` section of your `XFree86Config` file.

- The display engines of other operating systems.

The Unicode HOWTO

- The yudit editor, see below, and its printing engine.

Some no-cost TrueType fonts with large Unicode coverage are

Bitstream Cyberbit

Covers Roman, Cyrillic, Greek, Hebrew, Arabic, combining diacritical marks, Chinese, Korean, Japanese, and more.

Downloadable from <ftp://ftp.netscape.com/pub/communicator/extras/fonts/windows/Cyberbit.ZIP>. It is free for non-commercial purposes.

Microsoft Arial

Covers Roman, Cyrillic, Greek, Hebrew, Arabic, some combining diacritical marks, Vietnamese.

Downloadable; look on a search engine for ftp-able files called `arial.ttf`, `ariali.ttf`, `arialbd.ttf`, `arialbi.ttf`.

Lucida Sans Unicode

Covers Roman, Cyrillic, Greek, Hebrew, combining diacritical marks.

Download: contained in IBM's JDK 1.3.0 for Linux, at <http://www.ibm.com/java/jdk/linux130/>, or directly downloadable as `LucidaSansRegular.ttf` and `LucidaSansOblique.ttf` from <ftp://ftp.maths.tcd.ie/Linux/opt/IBMJava2-13/jre/lib/fonts/>.

Arphic

Cover Chinese (both traditional and simplified).

Download: at <ftp://ftp.gnu.org/non-gnu/chinese-fonts-truetype/>. These fonts are truly free.

Download locations for these and other TrueType fonts can be found at Christoph Singer's list of freely downloadable Unicode TrueType fonts <http://www.ccss.de/slovo/unifonts.htm>.

Truetype fonts are installed similarly to fixed size fonts, except that they go in a separate directory, and that `ttmkfdir` must be called before `mkfontdir`:

```
# mkdir -p /usr/X11R6/lib/X11/fonts/truetype
# cp /somewhere/Cyberbit.ttf ... /usr/X11R6/lib/X11/fonts/truetype
# cd /usr/X11R6/lib/X11/fonts/truetype
# ttmkfdir > fonts.scale
# mkfontdir
# xset fp rehash
```

TrueType fonts can be converted to low resolution, non-scalable X11 fonts by use of Mark Leisher's `tft2bdf` utility <ftp://crl.nmsu.edu/CLR/multiling/General/tft2bdf-2.8-LINUX.tar.gz>. For example, to generate a proportional Unicode font for use with `cooledit`:

```
# cd /usr/X11R6/lib/X11/fonts/local
# tft2bdf ../truetype/Cyberbit.ttf > cyberbit.bdf
# bdf2pcf -o cyberbit.pcf cyberbit.bdf
# gzip -9 cyberbit.pcf
# mkfontdir
# xset fp rehash
```

More information about TrueType fonts can be found in the Linux TrueType HOWTO <http://www.moisty.org/~brion/linux/TrueType-HOWTO.html>.

2.6 Miscellaneous

A small program which tests whether a Linux console or xterm is in UTF-8 mode can be found in the <ftp://sunsite.unc.edu/pub/Linux/system/keyboards/x-It-1.24.tar.gz> package by Ricardas Cepas, files testUTF-8.c and testUTF8.c. Most applications should not use this, however: they should look at the environment variables, see section "Locale environment variables".

3. Locale setup

3.1 Files & the kernel

You can now already use any Unicode characters in file names. No kernel or file utilities need modifications. This is because file names in the kernel can be anything not containing a null byte, and '/' is used to delimit subdirectories. When encoded using UTF-8, non-ASCII characters will never be encoded using null bytes or slashes. All that happens is that file and directory names occupy more bytes than they contain characters. For example, a filename consisting of five greek characters will appear to the kernel as a 10-byte filename. The kernel does not know (and does not need to know) that these bytes are displayed as greek.

This is the general theory, as long as your files stay inside Linux. On filesystems which are used from other operating systems, you have mount options to control conversion of filenames to/from UTF-8:

- The "vfat" filesystems has a mount option "utf8". See <file:/usr/src/linux/Documentation/filesystems/vfat.txt>. When you give an "iocharset" mount option different from the default (which is "iso8859-1"), the results with and without "utf8" are not consistent. Therefore I don't recommend the "iocharset" mount option.
- The "msdos", "umsdos" filesystems have the same mount option, but it appears to have no effect.
- The "iso9660" filesystem has a mount option "utf8". See <file:/usr/src/linux/Documentation/filesystems/isofs.txt>.
- Since Linux 2.2.x kernels, the "ntfs" filesystem has a mount option "utf8". See <file:/usr/src/linux/Documentation/filesystems/ntfs.txt>.

The other filesystems (nfs, smbfs, ncpfs, hpfs, etc.) don't convert filenames; therefore they support Unicode file names in UTF-8 encoding only if the other operating system supports them. Recall that to enable a mount option for all future remounts, you add it to the fourth column of the corresponding /etc/fstab line.

3.2 Upgrading the C library

glibc-2.2 supports multibyte locales, in particular UTF-8 locales. But glibc-2.1.x and earlier C libraries do not support it. Therefore you need to upgrade to glibc-2.2. Upgrading from glibc-2.1.x is riskless, because glibc-2.2 is binary compatible with glibc-2.1.x (at least on i386 platforms, and except for IPv6). Nevertheless, I recommend to have a bootable rescue disk handy in case something goes wrong.

Prepare the kernel sources. You must have them unpacked and configured. /usr/src/linux/include/linux/autoconf.h must exist. Building the kernel is not needed.

Retrieve the glibc sources <ftp://ftp.gnu.org/pub/gnu/glibc/>, su to root, then unpack, build and install it:

```
# unset LD_PRELOAD
# unset LD_LIBRARY_PATH
```

The Unicode HOWTO

```
# tar xvfz glibc-2.2.tar.gz
# tar xvfz glibc-linuxthreads-2.2.tar.gz -C glibc-2.2
# mkdir glibc-2.2-build
# cd glibc-2.2-build
# ../glibc-2.2/configure --prefix=/usr --with-headers=/usr/src/linux/include --enable-add-
# make
# make check
# make info
# LC_ALL=C make install
# make localedata/install-locales
```

Upgrading from glibc versions earlier than 2.1.x cannot be done this way; consider first installing a Linux distribution based on glibc-2.1.x, and then upgrading to glibc-2.2 as described above.

Note that if -- for any reason -- you want to rebuild GCC after having installed glibc-2.2, you need to first apply this patch [gcc-glibc-2.2-compat.diff](#) to the GCC sources.

3.3 General data conversion

You will need a program to convert your locally (probably ISO-8859-1) encoded texts to UTF-8. (The alternative would be to keep using texts in different encodings on the same machine; this is not fun in the long run.) One such program is `iconv`, which comes with glibc-2.2. Simply use

```
$ iconv --from-code=ISO-8859-1 --to-code=UTF-8 < old_file > new_file
```

Here are two handy shell scripts, called "i2u" [i2u.sh](#) (for ISO to UTF conversion) and "u2i" [u2i.sh](#) (for UTF to ISO conversion). Adapt according to your current 8-bit character set.

If you don't have glibc-2.2 and iconv installed, you can use GNU recode 3.6 instead. "i2u" [i2u_recode.sh](#) is "recode ISO-8859-1..UTF-8", and "u2i" [u2i_recode.sh](#) is "recode UTF-8..ISO-8859-1".
<ftp://ftp.gnu.org/pub/gnu/recode/recode-3.6.tar.gz>

Or you can also use CLISP instead. Here are "i2u" [i2u.lisp](#) and "u2i" [u2i.lisp](#) written in Lisp. Note: You need a CLISP version from July 1999 or newer. <ftp://clisp.cons.org/pub/lisp/clisp/source/clisp-src.tar.gz>.

Other data conversion programs, less powerful than GNU recode, are `trans`
<ftp://ftp.informatik.uni-erlangen.de/pub/doc/ISO/charsets/trans113.tar.gz>, `tcs` from the Plan9 operating system <ftp://ftp.informatik.uni-erlangen.de/pub/doc/ISO/charsets/tcs.tar.gz>, and `utrans`/`uhtrans`/`hutrans`
<ftp://ftp.cdrom.com/pub/FreeBSD/distfiles/i18ntools-1.0.tar.gz> by G. Adam Stanislav
<adam@whizkidtech.net>.

For the repeated conversion of files to UTF-8 from different character sets, a semi-automatic tool can be used: [to-utf8](#) presents the non-ASCII parts of a file to the user, lets him decide about the file's original character set, and then converts the file to UTF-8.

3.4 Locale environment variables

You may have the following environment variables set, containing locale names:

LANGUAGE

override for LC_MESSAGES, used by GNU gettext only

LC_ALL

override for all other LC_* variables

LC_CTYPE, LC_MESSAGES, LC_COLLATE, LC_NUMERIC, LC_MONETARY, LC_TIME

individual variables for: character types and encoding, natural language messages, sorting rules, number formatting, money amount formatting, date and time display

LANG

default value for all LC_* variables

(See `man 7 locale` for a detailed description.)

Each of the LC_* and LANG variables can contain a locale name of the following form:

```
language[_territory[.codeset]][@modifier]
```

where language is an [ISO 639](#) language code (lower case), territory is an [ISO 3166](#) country code (upper case), codeset denotes a character set, and modifier stands for other particular attributes (for example indicating a particular language dialect, or a nonstandard orthography).

LANGUAGE can contain several locale names, separated by colons.

In order to tell your system and all applications that you are using UTF-8, you need to add a codeset suffix of UTF-8 to your locale names. For example, if you were using

```
LC_CTYPE=de_DE
```

you would change this to

```
LC_CTYPE=de_DE.UTF-8
```

You do *not* need to change your LANGUAGE environment variable. GNU gettext in glibc-2.2 has the ability to convert translations to the right encoding.

3.5 Creating the locale support files

You create using `localedef` the support files for each UTF-8 locale you intend to use, for example:

```
$ localedef -v -c -i de_DE -f UTF-8 de_DE.UTF-8
```

You typically don't need to create locales named "de" or "fr" without country suffix, because these locales are normally only used by the LANGUAGE variable and not by the LC_* variables, and LANGUAGE is only used as an override for LC_MESSAGES.

4. Specific applications

4.1 Shells

bash

By default, GNU bash assumes that every character is one byte long and one column wide. A patch for bash 2.04, by Marcin 'Qrczak' Kowalczyk and Ricardas Cepas, teaches bash about multibyte characters in UTF-8

encoding. [bash-2.04-diff](#)

Double-width characters, combining characters and bidi are not supported by this patch. It seems a complete redesign of the readline redisplay engine is needed.

4.2 Networking

telnet

In some installations, telnet is not 8-bit clean by default. In order to be able to send Unicode keystrokes to the remote host, you need to set telnet into "outbinary" mode. There are two ways to do this:

```
$ telnet -L <host>
```

and

```
$ telnet
telnet> set outbinary
telnet> open <host>
```

kermit

The communications program C-Kermit <http://www.columbia.edu/kermit/ckermi.html>, (an interactive tool for connection setup, telnet, file transfer, with support for TCP/IP and serial lines), in versions 7.0 or newer, understands the file and transfer encodings UTF-8 and UCS-2, and understands the terminal encoding UTF-8, and converts between these encodings and many others. Documentation of these features can be found in <http://www.columbia.edu/kermit/ckermi2.html#x6.6>.

4.3 Browsers

Netscape

Netscape 4.05 or newer can display HTML documents in UTF-8 encoding. All a document needs is the following line between the <head> and </head> tags:

```
<meta http-equiv="Content-Type" content="text/html; charset=UTF-8">
```

Netscape 4.05 or newer can also display HTML and text files in UCS-2 encoding with byte-order mark.

<http://www.netscape.com/computing/download/>

Mozilla

Mozilla milestone M16 has much better internationalization than Netscape 4. It can display HTML documents in UTF-8 encoding with support for more languages. Alas, there is a cosmetic problem with CJK fonts: some glyphs can be bigger than the line's height, thus overlapping the previous or next line.

<http://www.mozilla.org/>

Amaya

Amaya 4.2.1 (<http://www.w3.org/Amaya/>, <http://www.w3.org/Amaya/User/SourceDist>) has now limited handling of UTF-8 encoded HTML pages. It recognizes the encoding, but it displays only ISO-8859-1 and symbol characters; it only ever accesses the fonts

```
-adobe-times-*-iso8859-1
-adobe-helvetica-*-iso8859-1
-adobe-new century schoolbook-*-iso8859-1
-adobe-courier-*-iso8859-1
-adobe-symbol-*-adobe-fontspecific
```

Amaya is in fact a HTML editor, not only a browser. Amaya's strengths among the browsers are its speed, given enough memory, and its rendering of mathematical formulas (MathML support).

lynx

lynx-2.8 has an options screen (key 'O') which permits to set the display character set. When running in an xterm or Linux console in UTF-8 mode, set this to "UNICODE UTF-8". Note that for this setting to take effect in the current browser session, you have to confirm on the "Accept Changes" field, and for this setting to take effect in future browser sessions, you have to enable the "Save options to disk" field and then confirm it on the "Accept Changes" field.

Now, again, all a document needs is the following line between the <head> and </head> tags:

```
<meta http-equiv="Content-Type" content="text/html; charset=UTF-8">
```

When you are viewing text files in UTF-8 encoding, you also need to pass the command-line option "-assume_local_charset=UTF-8" (affects only file:/... URLs) or "-assume_charset=UTF-8" (affects all URLs). In lynx-2.8.2 you can alternatively, in the options screen (key 'O'), change the assumed document character set to "utf-8".

There is also an option in the options screen, to set the "preferred document character set". But it has no effect, at least with file:/... URLs and with http://... URLs served by apache-1.3.0.

There is a spacing and line-breaking problem, however. (Look at the russian section of x-utf8.html, or at utf-8-demo.txt.)

Also, in lynx-2.8.2, configured with --enable-prettysrc, the nice colour scheme does not work correctly any more when the display character set has been set to "UNICODE UTF-8". This is fixed by a simple patch [lynx282.diff](#).

The Lynx developers say: "For any serious use of UTF-8 screen output with lynx, compiling with slang lib and -DSLING_MBCS_HACK is still recommended."

Latest stable release: <ftp://ftp.gnu.org/pub/gnu/lynx/lynx-2.8.2.tar.gz>

<http://lynx.isc.org/>

General home page: <http://lynx.browser.org/>

<http://www.slcc.edu/lynx/>

Newer development snapshots: <http://lynx.isc.org/current/>, <ftp://lynx.isc.org/current/>

w3m

w3m by Akinori Ito <http://ei5nazha.yz.yamagata-u.ac.jp/~aito/w3m/eng/> is a text mode browser for HTML pages and plain-text files. Its layout of HTML tables, enumerations etc. is much prettier than lynx' one. w3m can also be used as a high quality HTML to plain text converter.

w3m 0.1.10 has command line options for the three major Japanese encodings, but can also be used for UTF-8 encoded files. Without command line options, you often have to press Ctrl-L to refresh the display, and line breaking in Cyrillic and CJK paragraphs is not good.

To fix this, by Hironori Sakamoto has a patch <http://www2u.biglobe.ne.jp/~hsaka/w3m/> which adds UTF-8 as display encoding.

Test pages

Some test pages for browsers can be found at the pages of Alan Wood <http://www.hclrss.demon.co.uk/unicode/#links> and James Kass <http://home.att.net/~jameskass/>.

4.4 Editors

yudit

yudit by Gáspár Sinai <http://www.yudit.org/> is a first-class unicode text editor for the X Window System. It supports simultaneous processing of many languages, input methods, conversions for local character standards. It has facilities for entering text in all languages with only an English keyboard, using keyboard configuration maps.

yudit-1.5

It can be compiled in three versions: Xlib GUI, KDE GUI, or Motif GUI.

Customization is very easy. Typically you will first customize your font. From the font menu I chose "Unicode". Then, since the command "xlsfonts '*-*iso10646-1'" still showed some ambiguity, I chose a font size of 13 (to match Markus Kuhn's 13-pixel fixed font).

Next, you will customize your input method. The input methods "Straight", "Unicode" and "SGML" are most remarkable. For details about the other built-in input methods, look in /usr/local/share/yudit/data/.

To change the default for the next session, edit your \$HOME/.yuditrc file.

The general editor functionality is limited to editing, cut&paste and search&replace. No undo.

yudit-2.1

This version is less easy to learn, because it comes with a homebrewn GUI and no easily accessible help. But it has an undo functionality and should therefore be more usable than version 1.5.

Fonts for yudit

yudit can display text using a TrueType font; see section "TrueType fonts" above. The Bitstream Cyberbit gives good results. For yudit to find the font, symlink it to
`/usr/local/share/yudit/data/cyberbit.ttf`.

vim

vim (as of version 6.0r) has good support for UTF-8: when started in an UTF-8 locale, it assumes UTF-8 encoding for the console and the text files being edited. It supports double-wide (CJK) characters as well and combining characters and therefore fits perfectly into UTF-8 enabled xterm.

Installation: Download from <http://www.vim.org/>. After unpacking the four parts, call `./configure` with `--with-features=big --enable-multibyte` arguments (or edit `src/Makefile` to include the `--with-features=big` and `--enable-multibyte` options). This will turn on the feature `FEAT_MBYTE`. Then do "make" and "make install".

vim can be used to edit files in other encodings. For example, to edit a BIG5 encoded file: `:e ++cc=BIG5 filename`. All encoding names supported by `iconv` are accepted. Plus: vim automatically distinguishes UTF-8 and ISO-8859-1 files without needing any command line option.

cooledit

cooledit by Paul Sheer <http://www.cooledit.org/> is a good text editor for the X Window System. Since version 3.15, it has support for Unicode, including Bidi for Hebrew (but not Arabic).

A build error message message about a missing "vga_setpage" function is worked around by adding `"-DDO_NOT_USE_VGALIB"` to the `CFLAGS`.

To view UTF-8 files in an UTF-8 locale you have to modify a setting in the "Options -> Switches" panel: Enable the checkbox "Display characters outside locale". I also found it necessary to disable "Spellcheck as you type".

For viewing texts with both European and CJK characters, cooledit needs a font which contains both, for example the GNU unifont (see section "X11 Unicode fonts"): Start once

```
$ cooledit -fn -gnu-unifont-medium-r-normal--16-160-75-75-c-80-iso10646-1
```

cooledit will then use this font in all future invocations.

Unfortunately, the only characters that can be entered through the keyboard are ISO-8859-1 characters and, through a cooledit specific compose mechanism, ISO-8859-2 characters. Inputting arbitrary Unicode characters in cooledit is possible, but a bit tedious.

emacs

First of all, you should read the section "International Character Set Support" (node "International") in the Emacs manual. In particular, note that you need to start Emacs using the command

```
$ emacs -fn fontset-standard
```

so that it will use a font set comprising a lot of international characters.

In the short term, there are two packages for using UTF-8 in Emacs. None of them needs recompiling Emacs.

- The emacs-utf package <http://www.cs.ust.hk/faculty/otfried/Mule/> by Otfried Cheong provides a "unicode-utf8" encoding to Emacs.
- The oc-unicode package <http://www.cs.ust.hk/faculty/otfried/Mule/>, by Otfried Cheong, an extension of the Mule-UCS package <ftp://etlport.etl.go.jp/pub/mule/Mule-UCS/Mule-UCS-0.70.tar.gz> (mirrored at <http://riksun.riken.go.jp/archives/misc/mule/Mule-UCS/Mule-UCS-0.70.tar.gz> and <ftp://ftp.m17n.org/pub/mule/Mule-UCS/Mule-UCS-0.70.tar.gz>) by Hisashi Miyashita, provides a "utf-8" encoding to Emacs.

You can use either of these packages, or both together. The advantages of the emacs-utf "unicode-utf8" encoding are: it loads faster, and it deals better with combining characters (important for Thai). The advantage of the Mule-UCS / oc-unicode "utf-8" encoding is: it can apply to a process buffer (such as M-x shell), not only to loading and saving of files; and it respects the widths of characters better (important for Ethiopian). However, it is less reliable: After heavy editing of a file, I have seen some Unicode characters replaced with U+FFFD after the file was saved. (But maybe that were bugs in Emacs 20.5 and 20.6 which are fixed in Emacs 20.7.)

To install the emacs-utf package, compile the program "utf2mule" and install it somewhere in your \$PATH, also install unicode.el, muleuni-1.el, unicode-char.el somewhere. Then add the lines

```
(setq load-path (cons "/home/user/somewhere/emacs" load-path))
(if (not (string-match "XEmacs" emacs-version))
    (progn
      (require 'unicode)
      ;(setq unicode-data-path ".../UnicodeData-3.0.0.txt")
      (if (eq window-system 'x)
          (progn
            (setq fontset12
                  (create-fontset-from-fontset-spec
                    "-misc-fixed-medium-r-normal-*-12-*-*-*-*fontset-standard"))
            (setq fontset13
                  (create-fontset-from-fontset-spec
                    "-misc-fixed-medium-r-normal-*-13-*-*-*-*fontset-standard"))
            (setq fontset14
                  (create-fontset-from-fontset-spec
                    "-misc-fixed-medium-r-normal-*-14-*-*-*-*fontset-standard"))
            (setq fontset15
                  (create-fontset-from-fontset-spec
                    "-misc-fixed-medium-r-normal-*-15-*-*-*-*fontset-standard"))
            (setq fontset16
                  (create-fontset-from-fontset-spec
                    "-misc-fixed-medium-r-normal-*-16-*-*-*-*fontset-standard"))
            (setq fontset18
                  (create-fontset-from-fontset-spec
                    "-misc-fixed-medium-r-normal-*-18-*-*-*-*fontset-standard"))
```

The Unicode HOWTO

```
; (set-default-font fontset15)
)))))
```

to your `$HOME/.emacs` file. To activate any of the font sets, use the Mule menu item "Set Font/FontSet" or Shift-down-mouse-1. The Unicode coverage may of the font sets at different sizes may depend on the installed fonts; here are screen shots at various sizes of UTF-8-demo.txt ([12](#), [13](#), [14](#), [15](#), [16](#), [18](#)) and of the Mule script examples ([12](#), [13](#), [14](#), [15](#), [16](#), [18](#)). To designate a font set as the initial font set for the first frame at startup, uncomment the `set-default-font` line in the code snippet above.

To install the `oc-unicode` package, execute the command

```
$ emacs -batch -l oc-comp.el
```

and install the resulting file `un-define.elc`, as well as `oc-unicode.el`, `oc-charsets.el`, `oc-tools.el`, somewhere. Then add the lines

```
(setq load-path (cons "/home/user/somewhere/emacs" load-path))
(if (not (string-match "XEmacs" emacs-version))
    (progn
      (require 'oc-unicode)
      ;(setq unicode-data-path ".../UnicodeData-3.0.0.txt")
      (if (eq window-system 'x)
          (progn
            (setq fontset12
                  (oc-create-fontset
                   "-misc-fixed-medium-r-normal-*-12-*-*-*-*fontset-standard"
                   "-misc-fixed-medium-r-normal-ja-12-*-iso10646-*"))
            (setq fontset13
                  (oc-create-fontset
                   "-misc-fixed-medium-r-normal-*-13-*-*-*-*fontset-standard"
                   "-misc-fixed-medium-r-normal-ja-13-*-iso10646-*"))
            (setq fontset14
                  (oc-create-fontset
                   "-misc-fixed-medium-r-normal-*-14-*-*-*-*fontset-standard"
                   "-misc-fixed-medium-r-normal-ja-14-*-iso10646-*"))
            (setq fontset15
                  (oc-create-fontset
                   "-misc-fixed-medium-r-normal-*-15-*-*-*-*fontset-standard"
                   "-misc-fixed-medium-r-normal-ja-15-*-iso10646-*"))
            (setq fontset16
                  (oc-create-fontset
                   "-misc-fixed-medium-r-normal-*-16-*-*-*-*fontset-standard"
                   "-misc-fixed-medium-r-normal-ja-16-*-iso10646-*"))
            (setq fontset18
                  (oc-create-fontset
                   "-misc-fixed-medium-r-normal-*-18-*-*-*-*fontset-standard"
                   "-misc-fixed-medium-r-normal-ja-18-*-iso10646-*"))
            ; (set-default-font fontset15)
          )))
```

to your `$HOME/.emacs` file. You can choose your appropriate font set as with the `emacs-utf` package.

In order to open an UTF-8 encoded file, you will type

```
M-x universal-coding-system-argument unicode-utf8 RET
M-x find-file filename RET
```

or

The Unicode HOWTO

```
C-x RET c unicode-utf8 RET
C-x C-f filename RET
```

(or utf-8 instead of unicode-utf8, if you prefer oc-unicode/Mule-UCS).

In order to start a shell buffer with UTF-8 I/O, you will type

```
M-x universal-coding-system-argument utf-8 RET
M-x shell RET
```

(This works with oc-unicode/Mule-UCS only.)

There is a newer version Mule-UCS-0.81. Unfortunately you need to rebuild emacs from source in order to use it.

Note that all this works with Emacs 20 in windowing mode only, not in terminal mode. None of the mentioned packages works in Emacs 21, as of this writing.

Richard Stallman plans to add integrated UTF-8 support to Emacs in the long term, and so does the XEmacs developers group.

xemacs

(This section is written by Gilbert Baumann.)

Here is how to teach XEmacs (20.4 configured with MULE) the UTF-8 encoding. Unfortunately you need its sources to be able to patch it.

First you need these files provided by Tomohiko Morioka:

<http://turnbull.sk.tsukuba.ac.jp/Tools/XEmacs/xemacs-21.0-b55-emc-b55-ucs.diff> and
<http://turnbull.sk.tsukuba.ac.jp/Tools/XEmacs/xemacs-ucs-conv-0.1.tar.gz>

The .diff is a diff against the C sources. The tar ball is elisp code, which provides lots of code tables to map to and from Unicode. As the name of the diff file suggests it is against XEmacs-21; I needed to help `patch' a bit. The most notable difference to my XEmacs-20.4 sources is that file-coding.[ch] was called mule-coding.[ch].

For those unfamiliar with the XEmacs-MULE stuff (as I am) a quick guide:

What we call an encoding is called by MULE a `coding-system'. The most important commands are:

```
M-x set-file-coding-system
M-x set-buffer-process-coding-system [comint buffers]
```

and the variable `file-coding-system-alist', which guides `find-file' to guess the encoding used. After stuff was running, the very first thing I did was [this](#).

This code looks into the special mode line introduced by -*- somewhere in the first 600 bytes of the file about to opened; if now there is a field "Encoding: xyz;" and the xyz encoding ("coding system" in Emacs speak) exists, choose that. So now you could do e.g.

```
;;; -*- Mode: Lisp; Syntax: Common-Lisp; Package: CLEX; Encoding: utf-8; -*-
```

The Unicode HOWTO

and XEmacs goes into utf-8 mode here.

After everything was running I defined \u03BB (greek lambda) as a macro like:

```
(defmacro \u03BB (x) `(lambda .,x))
```

nedit

xedit

With XFree86-4.0.1, xedit is able to edit UTF-8 files if you set the locale accordingly (see above), and add the line "Xedit*international: true" to your \$HOME/.Xdefaults file.

axe

As of version 6.1.2, aXe supports only 8-bit locales. If you add the line "Axe*international: true" to your \$HOME/.Xdefaults file, it will simply dump core.

pico

As of version 4.30, pine cannot be reasonably used to view or edit UTF-8 files. In UTF-8 enabled xterm, it has severe redraw problems.

mined98

mined98 is a small text editor by Michiel Huisjes, Achim Müller and Thomas Wolff.

<http://www.inf.fu-berlin.de/~wolff/mined98.tar.gz> It lets you edit UTF-8 or 8-bit encoded files, in an UTF-8 or 8-bit xterm. It also has powerful capabilities for entering Unicode characters.

mined lets you edit both 8-bit encoded and UTF-8 encoded files. By default it uses an autodetection heuristic. If you don't want to rely on heuristics, pass the command-line option `-u` when editing an UTF-8 file, or `+u` when editing an 8-bit encoded file. You can change the interpretation at any time from within the editor: It displays the encoding ("L:h" for 8-bit, "U:h" for UTF-8) in the menu line. Click on the first of these characters to change it.

mined knows about double-width and combining characters and displays them correctly. It also has a special display mode for combining characters.

mined also has a scrollbar and very nice pull-down menus. Alas, the "Home", "End", "Delete" keys do not work.

qemacs

qemacs 0.2 is a small text editor by Fabrice Bellard. <http://www-stud.enst.fr/~bellard/qemacs/> with Emacs keybindings. It runs in an UTF-8 console or xterm, and can edit both 8-bit encoded and UTF-8 encoded files. It still has a few rough edges, but further development is underway.

4.5 Mailers

MIME: RFC 2279 defines UTF-8 as a MIME charset, which can be transported under the 8bit, quoted-printable and base64 encodings. The older MIME UTF-7 proposal (RFC 2152) is considered to be deprecated and should not be used any further.

Mail clients released after January 1, 1999, should be capable of sending and displaying UTF-8 encoded mails, otherwise they are considered deficient. But these mails have to carry the MIME labels

```
Content-Type: text/plain; charset=UTF-8
Content-Transfer-Encoding: 8bit
```

Simply piping an UTF-8 file into "mail" without caring about the MIME labels will not work.

Mail client implementors should take a look at <http://www.imc.org/imc-intl/> and <http://www.imc.org/mail-i18n.html>.

Now about the individual mail clients (or "mail user agents"):

pine

The situation for an unpatched pine version 4.30 is as follows.

Pine does not do character set conversions. But it allows you to view UTF-8 mails in an UTF-8 text window (Linux console or xterm).

Normally, Pine will warn about different character sets each time you view an UTF-8 encoded mail. To get rid of this warning, choose S (setup), then C (config), then change the value of "character-set" to UTF-8. This option will not do anything, except to reduce the warnings, as Pine has no built-in knowledge of UTF-8.

Also note that Pine's notion of Unicode characters is pretty limited: It will display Latin and Greek characters, but not other kinds of Unicode characters.

A patch by Robert Brady <robert@suse.co.uk> <http://www.ents.susu.soton.ac.uk/~robert/pine-utf8-0.1.diff> adds UTF-8 support to Pine. With this patch, it decodes and prints headers and bodies properly. The patch depends on the GNOME libunicode <http://cvs.gnome.org/lxr/source/libunicode/>.

However, alignment remains broken in many places; replying to a mail does not cause the character set to be converted as appropriate; and the editor, pico, cannot deal with multibyte characters.

kmail

kmail (as of KDE 1.0) does not support UTF-8 mails at all.

Netscape Communicator

Netscape Communicator's Messenger can send and display mails in UTF-8 encoding, but it needs a little bit of manual user intervention.

The Unicode HOWTO

To send an UTF-8 encoded mail: After opening the "Compose" window, but before starting to compose the message, select from the menu "View -> Character Set -> Unicode (UTF-8)". Then compose the message and send it.

When you receive an UTF-8 encoded mail, Netscape unfortunately does not display it in UTF-8 right away, and does not even give a visual clue that the mail was encoded in UTF-8. You have to manually select from the menu "View -> Character Set -> Unicode (UTF-8)".

For displaying UTF-8 mails, Netscape uses different fonts. You can adjust your font settings in the "Edit -> Preferences -> Fonts" dialog; choose the "Unicode" font category.

emacs (rmail, vm)

mutt

mutt-1.2.x, as available from <http://www.mutt.org/>, has only rudimentary support for UTF-8: it can convert from UTF-8 into an 8-bit display charset. The mutt-1.3.x development branch also supports UTF-8 as the display charset, so you can run Mutt in an UTF-8 xterm, and has thorough support for MIME and charset conversion (relying on iconv).

exmh

exmh 2.1.2 with Tk 8.4a1 can recognize and correctly display UTF-8 mails (without CJK characters) if you add the following lines to your `$HOME/.Xdefaults` file.

```
!
! Exmh
!
exmh.mimeUCharsets:           utf-8
exmh.mime_utf-8_registry:     iso10646
exmh.mime_utf-8_encoding:     1
exmh.mime_utf-8_plain_families: fixed
exmh.mime_utf-8_fixed_families: fixed
exmh.mime_utf-8_proportional_families: fixed
exmh.mime_utf-8_title_families: fixed
```

4.6 Text processing

groff

groff 1.16.1, the GNU implementation of the traditional Unix text processing system troff/nroff, can output UTF-8 formatted text. Simply use ``groff -Tutf8'` instead of ``groff -Tlatin1'` or ``groff -Tascii'`.

TeX

The teTeX 0.9 (and newer) distribution contains an Unicode adaptation of TeX, called Omega (<http://www.gutenberg.eu.org/omega/>, <ftp://ftp.ens.fr/pub/tex/yannis/omega>). Together with the `unicode.tex` file contained in [utf8-tex-0.1.tar.gz](http://www.gutenberg.eu.org/omega/) it enables you to use UTF-8 encoded sources as input for TeX. A thousand of Unicode characters are currently supported.

The Unicode HOWTO

All that changes is that you run ``omega'` (instead of ``tex'`) or ``lambda'` (instead of ``latex'`), and insert the following lines at the head of your source input.

```
\ocp\TexUTF=inutf8
\InputTranslation currentfile \TexUTF

\input unicode
```

Other maybe related links: <http://www.dante.de/projekte/nts/NTS-FAQ.html>,
<ftp://ftp.dante.de/pub/tex/language/chinese/CJK/>.

4.7 Databases

PostgreSQL

PostgreSQL 6.4 or newer can be built with the configuration option `--with-mb=UNICODE`.

Interbase

Borland/Inprise's Interbase 6.0 can store string fields in UTF-8 format if the option "CHARACTER SET UNICODE_FSS" is given.

4.8 Other text-mode applications

less

With <http://www.flash.net/~marknu/less/less-358.tar.gz> you can browse UTF-8 encoded text files in an UTF-8 xterm or console. Make sure that the environment variable LESSCHARSET is not set (or is set to utf-8). If you also have a LESSKEY environment variable set, also make sure that the file it points to does not define LESSCHARSET. If necessary, regenerate this file using the ``lesskey'` command, or unset the LESSKEY environment variable.

lv

lv-4.49.3 by Tomio Narita <http://www.ff.ij4u.or.jp/~nrt/lv/> is a file viewer with builtin character set converters. To view UTF-8 files in an UTF-8 console, use "lv -Au8". But it can also be used to view files in other CJK encodings in an UTF-8 console.

There is a small glitch: lv turns off xterm's cursor and doesn't turn it on again.

expand

Get the GNU textutils-2.0 and apply the patch [textutils-2.0.diff](#), then configure, add "#define HAVE_FGETWC 1", "#define HAVE_FPUTWC 1" to config.h. Then rebuild.

col, colcrt, colrm, column, rev, ul

Get the util-linux-2.9y package, configure it, then define ENABLE_WIDECHAR in defines.h, change the "#if 0" to "#if 1" in lib/widechar.h. In text-utils/Makefile, modify CFLAGS and LDFLAGS so that they include

the directories where libutf8 is installed. Then rebuild.

figlet

figlet 2.2 has an option for UTF-8 input: "figlet -C utf8"

Base utilities

The Linux list of commands and utilities that ought to be made interoperable with UTF-8 is as follows. Useful information needs to get added here; I just didn't get around it yet :-)

As of glibc-2.2, regular expressions only work for 8-bit characters. In an UTF-8 locale, regular expressions that contain non-ASCII characters or that expect to match a single multibyte character with "." do not work. This affects all commands and utilities listed below.

alias

No info available yet.

ar

No info available yet.

arch

No info available yet.

arp

No info available yet.

at

As of at-3.1.8: The two uses of isalnum in at.c are invalid and should be replaced with a use of quotearg.c or an exclude list of the (fixed) list of shell metacharacters. The two uses of %8s in at.c and atd.c are invalid and should become arbitrary length.

awk

No info available yet.

basename

As of sh-utils-2.0i: OK.

batch

No info available yet.

bc

No info available yet.

bg

No info available yet.

bunzip2

No info available yet.

bzip2

No info available yet.

bzip2recover

No info available yet.

cal

No info available yet.

cat

No info available yet.

cd

No info available yet.

cflow

	No info available yet.
chgrp	As of fileutils-4.0u: OK.
chmod	As of fileutils-4.0u: OK.
chown	As of fileutils-4.0u: OK.
chroot	As of sh-utils-2.0i: OK.
cksum	As of textutils-2.0e: OK.
clear	No info available yet.
cmp	No info available yet.
col	No info available yet.
comm	No info available yet.
command	No info available yet.
compress	No info available yet.
cp	As of fileutils-4.0u: OK.
cpio	No info available yet.
crontab	No info available yet.
csplit	No info available yet.
ctags	No info available yet.
cut	No info available yet.
date	As of sh-utils-2.0i: OK.
dd	As of fileutils-4.0u: The conv=lcase, conv=ucase options don't work correctly.
df	As of fileutils-4.0u: OK.
diff	As of diffutils-2.7.2: the --side-by-side mode therefore doesn't compute column width correctly.
diff3	No info available yet.
dirname	As of sh-utils-2.0i: OK.
domainname	No info available yet.
du	As of fileutils-4.0u: OK.

echo	As of sh-utils-2.0i: OK.
ed	No info available yet.
egrep	No info available yet.
env	As of sh-utils-2.0i: OK.
ex	No info available yet.
expand	No info available yet.
expr	As of sh-utils-2.0i: The operators "match", "substr", "index", "length" don't work correctly.
false	As of sh-utils-2.0i: OK.
fc	No info available yet.
fg	No info available yet.
fgrep	No info available yet.
file	No info available yet.
find	As of findutils-4.1.6: The "-iregex" does not work correctly; this needs a fix in function find/parser.c:insert_regex.
fold	No info available yet.
ftp[BSD]	No info available yet.
fuser	No info available yet.
gencat	No info available yet.
getconf	No info available yet.
getopts	No info available yet.
gettext	No info available yet.
grep	No info available yet.
gunzip	No info available yet.
gzip	gzip-1.3 is UTF-8 capable, but it uses only English messages in ASCII charset. Proper internationalization would require: Use gettext. Call setlocale. In function check_ofname (file gzip.c), use the function rpmatch from GNU text/sh/fileutils instead of asking for "y" or "n". The use of strlen in gzip.c:852 is wrong, needs to use the function mbswidth.
hash	

No info available yet.
head
No info available yet.
hostname
As of sh-utils-2.0i: OK.
iconv
No info available yet.
id
As of sh-utils-2.0i: OK.
ifconfig
No info available yet.
imake
No info available yet.
ipcrm
No info available yet.
ipcs
No info available yet.
jobs
No info available yet.
join
No info available yet.
kill
No info available yet.
killall
No info available yet.
ldd
No info available yet.
less
No complete info available yet.
lex
No info available yet.
ln
As of fileutils-4.0u: OK.
locale
As of glibc-2.2: OK.
localedef
As of glibc-2.2: OK.
logger
No info available yet.
logname
As of sh-utils-2.0i: OK.
lp
No info available yet.
lpc[BSD]
No info available yet.
lpq[BSD]
No info available yet.
lpr[BSD]
No info available yet.
lprm[BSD]
No info available yet.

lpstat(LEGACY)

No info available yet.

ls

As of fileutils-4.0y: OK.

m4

No info available yet.

mailx

No info available yet.

make

No info available yet.

man

No info available yet.

mesg

No info available yet.

mkdir

As of fileutils-4.0u: OK.

mkfifo

As of fileutils-4.0u: OK.

mkfs

No info available yet.

mkswap

No info available yet.

more

No info available yet.

mount

No info available yet.

msgfmt

No info available yet.

msgmerge

No info available yet.

mv

As of fileutils-4.0u: OK.

netstat

No info available yet.

newgrp

No info available yet.

nice

As of sh-utils-2.0i: OK.

nl

No info available yet.

nohup

As of sh-utils-2.0i: OK.

nslookup

No info available yet.

nm

No info available yet.

od

No info available yet.

passwd[BSD]

No info available yet.

paste

No info available yet.
patch
No info available yet.
pathchk
As of sh-utils-2.0i: OK.
ping
No info available yet.
pr
No info available yet.
printf
As of sh-utils-2.0i: OK.
ps
No info available yet.
pwd
As of sh-utils-2.0i: OK.
read
No info available yet.
reboot
No info available yet.
renice
No info available yet.
rm
As of fileutils-4.0u: OK.
rmdir
As of fileutils-4.0u: OK.
sed
No info available yet.
shar[BSD]
No info available yet.
shutdown
No info available yet.
sleep
As of sh-utils-2.0i: OK.
sort
No info available yet.
split
No info available yet.
strings
No info available yet.
strip
No info available yet.
stty
As of sh-utils-2.0.11: OK.
su[BSD]
No info available yet.
sum
As of textutils-2.0e: OK.
tail
No info available yet.
talk
No info available yet.

tar As of tar-1.13.17: OK, if user and group names are always ASCII.

tcsh No info available yet.

tee As of sh-utils-2.0i: OK.

telnet No info available yet.

test As of sh-utils-2.0i: OK.

time No info available yet.

touch As of fileutils-4.0u: OK.

tput No info available yet.

tr No info available yet.

true As of sh-utils-2.0i: OK.

tsort No info available yet.

tty As of sh-utils-2.0i: OK.

type No info available yet.

ulimit No info available yet.

umask No info available yet.

umount No info available yet.

unalias No info available yet.

uname As of sh-utils-2.0i: OK.

uncompress No info available yet.

unexpand No info available yet.

uniq No info available yet.

uudecode No info available yet.

uuencode No info available yet.

vi No info available yet.

wait No info available yet.

wc No info available yet.

	As of textutils-2.0.8: OK.
who	As of sh-utils-2.0i: OK.
wish	No info available yet.
write	No info available yet.
xargs	As of findutils-4.1.5: The program uses strstr; a patch has been submitted to the maintainer.
xgettext	No info available yet.
yacc	No info available yet.
zcat	No info available yet.

4.9 Other X11 applications

Owen Taylor is currently developing a library for rendering multilingual text, called pango.
<http://www.labs.redhat.com/~otaylor/pango/>, <http://www.pango.org/>.

5. Printing

Since Postscript itself does not support Unicode fonts, the burden of Unicode support in printing is on the program creating the Postscript document, not on the Postscript renderer.

The existing Postscript fonts I've seen - .pfa/.pfb/.afm/.pfm/.gsf - support only a small range of glyphs and are not Unicode fonts.

5.1 Printing using TrueType fonts

Both the uniprint and wprint programs produce good printed output for Unicode plain text. They require a TrueType font; see section "TrueType fonts" above. The Bitstream Cyberbit gives good results.

uniprint

The "uniprint" program contained in the yudit package can convert a text file to Postscript. For uniprint to find the Cyberbit font, symlink it to `/usr/local/share/yudit/data/cyberbit.ttf`.

wprint

The "wprint" (WorldPrint) program by Eduardo Trapani
<http://ttt.esperanto.org.uy/programoj/angle/wprint.html> postprocesses Postscript output produced by Netscape Communicator or Mozilla from HTML pages or plain text files.

The output is nearly perfect; only in Cyrillic paragraphs the line breaking is incorrect: the lines are only about half as wide as they should be.

Comparison

For plain text, uniprint has a better overall layout. On the other hand, only wprint gets Thai output correct.

5.2 Printing using fixed-size fonts

Generally, printing using fixed-size fonts does not give an as professional output as using TrueType fonts.

txtbdf2ps

The txtbdf2ps 0.7 program by Serge Winitzki <http://members.linuxstart.com/~winitzki/txtbdf2ps.html> converts a plain text file to Postscript, by use of a BDF font. Installation:

```
# install -m 777 txtbdf2ps-dev.txt /usr/local/bin/txtbdf2ps
```

Example with a proportional font:

```
$ txtbdf2ps -BDF=cyberbit.bdf -UTF-8 -nowrap < input.txt > output.ps
```

Example with a fixed-width font:

```
$ txtbdf2ps -BDF=unifont.bdf -UTF-8 -nowrap < input.txt > output.ps
```

Note: txtbdf2ps does not support combining characters and bidi.

5.3 The classical approach

Another way to print with TrueType fonts is to convert the TrueType font to a Postscript font using the `ttf2pt1` utility (<http://www.netSPACE.net.au/~mheath/ttf2pt1/>, <http://quadrant.netSPACE.net.au/ttf2pt1/>, <http://ttf2pt1.sourceforge.net/>). Details can be found in Julius Chroboczek's "Printing with TrueType fonts in Unix" writeup, <http://www.dcs.ed.ac.uk/home/jec/programs/xfsft/printing.html>.

TeX, Omega

TODO: CJK, metafont, omega, dvips, odvips, utf8-tex-0.1

DocBook

TODO: db2ps, jadetex

groff -Tps

"groff -Tps" produces Postscript output. Its Postscript output driver supports only a very limited number of Unicode characters (only what Postscript supports by itself).

5.4 No luck with...

Netscape's "Print..."

As of version 4.72, Netscape Communicator cannot correctly print HTML pages in UTF-8 encoding. You really have to use wprint.

Mozilla's "Print..."

As of version M16, printing of HTML pages is apparently not implemented.

html2ps

As of version 1.0b1, the html2ps HTML to Postscript converter does not support UTF-8 encoded HTML pages and has no special treatment of fonts: the generated Postscript uses the standard Postscript fonts.

a2ps

As of version 4.12, a2ps doesn't support printing UTF-8 encoded text.

enscript

As of version 1.6.1, enscript doesn't support printing UTF-8 encoded text. By default, it uses only the standard Postscript fonts, but it can also include a custom Postscript font in the output.

6. Making your programs Unicode aware

6.1 C/C++

The C `char` type is 8-bit and will stay 8-bit because it denotes the smallest addressable data unit. Various facilities are available:

For normal text handling

The ISO/ANSI C standard contains, in an amendment which was added in 1995, a "wide character" type `wchar_t`, a set of functions like those found in `<string.h>` and `<ctype.h>` (declared in `<wchar.h>` and `<wctype.h>`, respectively), and a set of conversion functions between `char *` and `wchar_t *` (declared in `<stdlib.h>`).

Good references for this API are

- the GNU libc-2.1 manual, chapters 4 "Character Handling" and 6 "Character Set Handling",
- the manual pages [man-mbswcs.tar.gz](http://ftp.win.tue.nl/pub/linux-local/manpages/man-pages-1.29.tar.gz), now contained in [ftp://ftp.win.tue.nl/pub/linux-local/manpages/man-pages-1.29.tar.gz](http://ftp.win.tue.nl/pub/linux-local/manpages/man-pages-1.29.tar.gz),
- the OpenGroup's introduction http://www.unix-systems.org/version2/whatsnew/login_mse.html,
- the OpenGroup's Single Unix specification <http://www.UNIX-systems.org/online.html>,
- the ISO/IEC 9899:1999 (ISO C 99) standard. The latest draft before it was adopted is called n2794. You find it at [ftp://ftp.csn.net/DMK/sc22wg14/review/](http://ftp.csn.net/DMK/sc22wg14/review/) or <http://java-tutor.com/docs/c/>.
- Clive Feather's introduction <http://www.lysator.liu.se/c/na1.html>,
- the Dinkumware C library reference http://www.dinkumware.com/htm_cl/.

The Unicode HOWTO

Advantages of using this API:

- It's a vendor independent standard.
- The functions do the right thing, depending on the user's locale. All a program needs to call is `setlocale(LC_ALL, "")` ;.

Drawbacks of this API:

- Some of the functions are not multithread-safe, because they keep a hidden internal state between function calls.
- There is no first-class locale datatype. Therefore this API cannot reasonably be used for anything that needs more than one locale or character set at the same time.
- The OS support for this API is not good on most OSes.

Portability notes

A `wchar_t` may or may not be encoded in Unicode; this is platform and sometimes also locale dependent. A multibyte sequence `char *` may or may not be encoded in UTF-8; this is platform and sometimes also locale dependent.

In detail, here is what the [Single Unix specification](#) says about the `wchar_t` type: *All wide-character codes in a given process consist of an equal number of bits. This is in contrast to characters, which can consist of a variable number of bytes. The byte or byte sequence that represents a character can also be represented as a wide-character code. Wide-character codes thus provide a uniform size for manipulating text data. A wide-character code having all bits zero is the null wide-character code, and terminates wide-character strings. The wide-character value for each member of the Portable Character Set (i.e. ASCII) will equal its value when used as the lone character in an integer character constant. Wide-character codes for other characters are locale- and implementation-dependent. State shift bytes do not have a wide-character code representation.*

One particular consequence is that in portable programs you shouldn't use non-ASCII characters in string literals. That means, even though you know the Unicode double quotation marks have the codes U+201C and U+201D, you shouldn't write a string literal `L"\u201cHello\u201d, he said"` or `"\xe2\x80\x9cHello\xe2\x80\x9d, he said"` in C programs. Instead, use GNU `gettext`, write it as `gettext("'Hello', he said")`, and create a message database `en.po` which translates "'Hello', he said" to `"\u201cHello\u201d, he said"`.

Here is a survey of the portability of the ISO/ANSI C facilities on various Unix flavours.

GNU `glibc-2.2.x`

- ◇ `<wchar.h>` and `<wctype.h>` exist.
- ◇ Has `wcs/mbs` functions, `fgetwc/fputwc/wprintf`, everything.
- ◇ Has five UTF-8 locales.
- ◇ `mbtowc` works.

GNU `glibc-2.0.x`, `glibc-2.1.x`

- ◇ `<wchar.h>` and `<wctype.h>` exist.
- ◇ Has `wcs/mbs` functions, but no `fgetwc/fputwc/wprintf`.
- ◇ No UTF-8 locale.

The Unicode HOWTO

◇ mbrtowc returns EILSEQ for bytes $\geq 0x80$.

AIX 4.3

- ◇ <wchar.h> and <wctype.h> exist.
- ◇ Has wcs/mbs functions, fgetwc/fputwc/wprintf, everything.
- ◇ Has many UTF-8 locales, one for every country.
- ◇ Needs `-D_XOPEN_SOURCE=500` in order to define `mbstate_t`.
- ◇ mbrtowc works.

Solaris 2.7

- ◇ <wchar.h> and <wctype.h> exist.
- ◇ Has wcs/mbs functions, fgetwc/fputwc/wprintf, everything.
- ◇ Has the following UTF-8 locales: `en_US.UTF-8`, `de.UTF-8`, `es.UTF-8`, `fr.UTF-8`, `it.UTF-8`, `sv.UTF-8`.
- ◇ mbrtowc returns `-1/EILSEQ` (instead of `-2`) for bytes $\geq 0x80$.

OSF/1 4.0d

- ◇ <wchar.h> and <wctype.h> exist.
- ◇ Has wcs/mbs functions, fgetwc/fputwc/wprintf, everything.
- ◇ Has an add-on `universal.utf8@ucs4` locale, see "man 5 unicode".
- ◇ mbrtowc does not know about UTF-8.

Irix 6.5

- ◇ <wchar.h> and <wctype.h> exist.
- ◇ Has wcs/mbs functions and fgetwc/fputwc, but not wprintf.
- ◇ Has no multibyte locales.
- ◇ Has only a dummy definition for `mbstate_t`.
- ◇ Doesn't have mbrtowc.

HP-UX 11.00

- ◇ <wchar.h> exists, <wctype.h> does not.
- ◇ Has wcs/mbs functions and fgetwc/fputwc, but not wprintf.
- ◇ Has a `C.utf8` locale.
- ◇ Doesn't have `mbstate_t`.
- ◇ Doesn't have mbrtowc.

As a consequence, I recommend to use the restartable and multithread-safe `wcsr/mbsr` functions, forget about those systems which don't have them (Irix, HP-UX, AIX), and use the UTF-8 locale plug-in `libutf8_plug.so` (see below) on those systems which permit you to compile programs which use these `wcsr/mbsr` functions (Linux, Solaris, OSF/1).

A similar advice, given by Sun in <http://www.sun.com/software/white-papers/wp-unicode/>, section "Internationalized Applications with Unicode", is:

To properly internationalize an application, use the following guidelines:

1. *Avoid direct access with Unicode. This is a task of the platform's internationalization framework.*
2. *Use the POSIX model for multibyte and wide-character interfaces.*
3. *Only call the APIs that the internationalization framework provides for language and cultural-specific operations.*
4. *Remain code-set independent.*

The Unicode HOWTO

If, for some reason, in some piece of code, you really have to assume that `wchar_t` is Unicode (for example, if you want to do special treatment of some Unicode characters), you should make that piece of code conditional upon the result of `is_locale_utf8()`. Otherwise you will mess up your program's behaviour in different locales or other platforms. The function `is_locale_utf8` is declared in [utf8locale.h](#) and defined in [utf8locale.c](#).

The libutf8 library

A portable implementation of the ISO/ANSI C API, which supports 8-bit locales and UTF-8 locales, can be found in [libutf8-0.7.3.tar.gz](#).

Advantages:

- Unicode UTF-8 support now, portably, even on OSes whose multibyte character support does not work or which don't have multibyte/wide character support at all.
- The same binary works in all OS supported 8-bit locales and in UTF-8 locales.
- When an OS vendor adds proper multibyte character support, you can take advantage of it by simply recompiling without `-DHAVE_LIBUTF8` compiler option.

The Plan9 way

The Plan9 operating system, a variant of Unix, uses UTF-8 as character encoding in all applications. Its wide character type is called `Rune`, not `wchar_t`. Parts of its libraries, written by Rob Pike and Howard Trickey, are available at <ftp://ftp.cdrom.com/pub/netlib/research/9libs/9libs-1.0.tar.gz>. Another similar library, written by Alistair G. Crooks, is <ftp://ftp.cdrom.com/pub/NetBSD/packages/distfiles/libutf-2.10.tar.gz>. In particular, each of these libraries contains an UTF-8 aware regular expression matcher.

Drawback of this API:

- UTF-8 is compiled in, not optional. Programs compiled in this universe lose support for the 8-bit encodings which are still frequently used in Europe.

For graphical user interface

The Qt-2.0 library <http://www.troll.no/> contains a fully-Unicode QString class. You can use the member functions `QString::utf8` and `QString::fromUtf8` to convert to/from UTF-8 encoded text. The `QString::ascii` and `QString::latin1` member functions should not be used any more.

For advanced text handling

The previously mentioned libraries implement Unicode aware versions of the ASCII concepts. Here are libraries which deal with Unicode concepts, such as titlecase (a third letter case, different from uppercase and lowercase), distinction between punctuation and symbols, canonical decomposition, combining classes, canonical ordering and the like.

ucdata-2.4

The `ucdata` library by Mark Leisher <http://crl.nmsu.edu/~mleisher/ucdata.html> deals with character properties, case conversion, decomposition, combining classes. The companion package `ure-0.5` <http://crl.nmsu.edu/~mleisher/ure-0.5.tar.gz> is a Unicode regular expression matcher.

ustring

The Unicode HOWTO

The `ustring` C++ library by Rodrigo Reyes <http://ustring.charabia.net/> deals with character properties, case conversion, decomposition, combining classes, and includes a Unicode regular expression matcher.

ICU

International Components for Unicode <http://oss.software.ibm.com/icu/>. IBM's very comprehensive internationalization library featuring Unicode strings, resource bundles, number formatters, date/time formatters, message formatters, collation and more. Lots of supported locales. Portable to Unix and Win32, but compiles out of the box only on Linux libc6, not libc5.

libunicode

The GNOME libunicode library <http://cvs.gnome.org/lxr/source/libunicode/> by Tom Tromeey and others. It covers character set conversion, character properties, decomposition.

For conversion

Two kinds of conversion libraries, which support UTF-8 and a large number of 8-bit character sets, are available:

iconv

The `iconv` implementation by Ulrich Drepper, contained in the GNU `glibc-2.2`.
<ftp://ftp.gnu.org/pub/gnu/glibc/glibc-2.2.tar.gz>. The `iconv` manpages are now contained in
<ftp://ftp.win.tue.nl/pub/linux-local/manpages/man-pages-1.29.tar.gz>.

The portable `iconv` implementation by Bruno Haible.
<ftp://ftp.ilog.fr/pub/Users/haible/gnu/libiconv-1.5.1.tar.gz>

The portable `iconv` implementation by Konstantin Chuguev. <joy@urc.ac.ru>
<ftp://ftp.urc.ac.ru/pub/local/OS/Unix/converters/iconv-0.4.tar.gz>

Advantages:

- `iconv` is POSIX standardized, programs using `iconv` to convert from/to UTF-8 will also run under Solaris. However, the names for the character sets differ between platforms. For example, "EUC-JP" under `glibc` is "eucJP" under HP-UX. (The official IANA name for this character set is "EUC-JP", so it's clearly a HP-UX deficiency.)
- On `glibc-2.1` systems, no additional library is needed. On other systems, one of the two other `iconv` implementations can be used.

librecode

`librecode` by François Pinard <ftp://ftp.gnu.org/pub/gnu/recode/recode-3.6.tar.gz>.

Advantages:

- Support for transliteration, i.e. conversion of non-ASCII characters to sequences of ASCII characters in order to preserve readability by humans, even when a lossless transformation is impossible.

Drawbacks:

- Non-standard API.

- Slow initialization.

ICU

International Components for Unicode 1.7 <http://oss.software.ibm.com/icu/>. IBM's internationalization library also has conversion facilities, declared in ``ucnv.h'`.

Advantages:

- Comprehensive set of supported encodings.

Drawbacks:

- Non-standard API.

Other approaches

libutf-8

libutf-8 by G. Adam Stanislav <adam@whizkidtech.net> contains a few functions for on-the-fly conversion from/to UTF-8 encoded ``FILE*'` streams.
<http://www.whizkidtech.net/i18n/libutf-8-1.0.tar.gz>

Advantages:

- ◇ Very small.

Drawbacks:

- ◇ Non-standard API.
- ◇ UTF-8 is compiled in, not optional. Programs compiled with this library lose support for the 8-bit encodings which are still frequently used in Europe.
- ◇ Installation is nontrivial: Makefile needs tweaking, not autoconfiguring.

6.2 Java

Java has Unicode support built into the language. The type ``char'` denotes a Unicode character, and the ``java.lang.String'` class denotes a string built up from Unicode characters.

Java can display any Unicode characters through its windowing system AWT, provided that 1. you set the Java system property "user.language" appropriately, 2. the `/usr/lib/java/lib/font.properties.language` font set definitions are appropriate, and 3. the fonts specified in that file are installed. For example, in order to display text containing japanese characters, you would install japanese fonts and run `"java -Duser.language=ja ..."`. You can combine font sets: In order to display western european, greek and japanese characters simultaneously, you would create a combination of the files "font.properties" (covers ISO-8859-1), "font.properties.el" (covers ISO-8859-7) and "font.properties.ja" into a single file. ??This is untested??

The interfaces `java.io.DataInput` and `java.io.DataOutput` have methods called ``readUTF'` and ``writeUTF'` respectively. But note that they don't use UTF-8; they use a modified UTF-8 encoding: the NUL character is encoded as the two-byte sequence `0xC0 0x80` instead of `0x00`, and a `0x00` byte is added at the end. Encoded this way, strings can contain NUL characters and nevertheless need not be prefixed with a length field - the C `<string.h>` functions like `strlen()` and `strcpy()` can be used to manipulate them.

6.3 Lisp

The Common Lisp standard specifies two character types: ``base-char'` and ``character'`. It's up to the implementation to support Unicode or not. The language also specifies a keyword argument ``:external-format'` to ``open'`, as the natural place to specify a character set or encoding.

Among the free Common Lisp implementations, only CLISP <http://clisp.cons.org/> supports Unicode. You need a CLISP version from March 2000 or newer. <ftp://clisp.cons.org/pub/lisp/clisp/source/clisp-src.tar.gz>. The types ``base-char'` and ``character'` are both equivalent to 16-bit Unicode. The functions `char-width` and `string-width` provide an API comparable to `wcwidth()` and `wcswidth()`. The encoding used for file or socket/pipe I/O can be specified through the ``:external-format'` argument. The encodings used for tty I/O and the default encoding for file/socket/pipe I/O are locale dependent.

Among the commercial Common Lisp implementations:

LispWorks http://www.xanalys.com/software_tools/products/ supports Unicode. The type ``base-char'` is equivalent to ISO-8859-1, and the type ``simple-char'` (subtype of ``character'`) contains all Unicode characters. The encoding used for file I/O can be specified through the ``:external-format'` argument, for example `(':UTF-8)`. Limitations: Encodings cannot be used for socket I/O. The editor cannot edit UTF-8 encoded files.

Eclipse <http://www.elwood.com/eclipse/eclipse.htm> supports Unicode. See <http://www.elwood.com/eclipse/char.htm>. The type ``base-char'` is equivalent to ISO-8859-1, and the type ``character'` contains all Unicode characters. The encoding used for file I/O can be specified through a combination of the ``:element-type'` and ``:external-format'` arguments to ``open'`. Limitations: Character attribute functions are locale dependent. Source and compiled source files cannot contain Unicode string literals.

The commercial Common Lisp implementation Allegro CL, in version 6.0, has Unicode support. The types ``base-char'` and ``character'` are both equivalent to 16-bit Unicode. The encoding used for file I/O can be specified through the ``:external-format'` argument, for example `:external-format :utf8`. The default encoding is locale dependent. More details are at <http://www.franz.com/support/documentation/6.0/doc/iacl.htm>.

6.4 Ada95

Ada95 was designed for Unicode support and the Ada95 standard library features special ISO 10646-1 data types `Wide_Character` and `Wide_String`, as well as numerous associated procedures and functions. The GNU Ada95 compiler (gnat-3.11 or newer) supports UTF-8 as the external encoding of wide characters. This allows you to use UTF-8 in both source code and application I/O. To activate it in the application, use `"WCEM=8"` in the FORM string when opening a file, and use compiler option `"-gnatW8"` if the source code is in UTF-8. See the GNAT (<ftp://cs.nyu.edu/pub/gnat/>) and Ada95 (<ftp://ftp.cnam.fr/pub/Ada/PAL/userdocs/docadalt/rm95/index.htm>) reference manuals for details.

6.5 Python

Python 2.0 (<http://www.python.org/2.0/>, <http://www.python.org/pipermail/python-announce-list/2000-October/000889.html>, <http://starship.python.net/crew/amk/python/writing/new-python/new-python.html>) contains Unicode support. It has a new fundamental data type ``unicode'`, representing a Unicode string, a module ``unicodedata'` for the character properties, and a set of converters for the most important encodings. See

<http://starship.python.net/crew/lemburg/unicode-proposal.txt>, or the file `Misc/unicode.txt` in the distribution, for details.

6.6 JavaScript/ECMAScript

Since JavaScript version 1.3, strings are always Unicode. There is no character type, but you can use the `\uXXXX` notation for Unicode characters inside strings. No normalization is done internally, so it expects to receive Unicode Normalization Form C, which the W3C recommends. See <http://developer.netscape.com/docs/manuals/communicator/jsref/js13.html#Unicode> for details and <http://developer.netscape.com/docs/javascript/e262-pdf.pdf> for the complete ECMAScript specification.

6.7 Tcl

Tcl/Tk started using Unicode as its base character set with version 8.1. Its internal representation for strings is UTF-8. It supports the `\uXXXX` notation for Unicode characters. See <http://dev.scriptics.com/doc/howto/i18n.html>.

6.8 Perl

Perl 5.6 stores strings internally in UTF-8 format, if you write

```
use utf8;
```

at the beginning of your script. `length()` returns the number of characters of a string. For details, see the Perl-i18n FAQ at <http://rf.net/~james/perli18n.html>.

Support for other (non-8-bit) encodings is available through the `iconv` interface module <http://cpan.perl.org/modules/by-module/Text/Text-Iconv-1.1.tar.gz>.

6.9 Related reading

Tomohiro Kubota has written an introduction to internationalization <http://www.debian.org/doc/manuals/intro-i18n/>. The emphasis of his document is on writing software that runs in any locale, using the locale's encoding.

7. Other sources of information

7.1 Mailing lists

Broader audiences can be reached at the following mailing lists.

Note that where I write ``at'`, you should write ``@'`. (Anti-spam device.)

linux-utf8

Address: `linux-utf8` at `nl.linux.org`

The Unicode HOWTO

This mailing list is about internationalization with Unicode, and covers a broad range of topics from the keyboard driver to the X11 fonts.

Archives are at <http://mail.nl.linux.org/linux-utf8/>.

To subscribe, send a message to `majordomo at nl.linux.org` with the line "subscribe linux-utf8" in the body.

li18nux

Address: `linux-i18n at sun.com`

This mailing list is focused on organizing internationalization work on Linux, and arranging meetings between people.

To subscribe, fill in the form at <http://www.li18nux.org/> and send it to `linux-i18n-request at sun.com`.

unicode

Address: `unicode at unicode.org`

This mailing list is focused on the standardization and continuing development of the Unicode standard, and related technologies, such as Bidi and sorting algorithms.

Archives are at <ftp://ftp.unicode.org/Public/MailArchive/>, but they are not regularly updated.

For subscription information, see <http://www.unicode.org/unicode/consortium/distlist.html>.

X11 internationalization

Address: `i18n at xfree86.org`

This mailing list addresses the people who work on better internationalization of the X11/XFree86 system.

Archives are at <http://devel.xfree86.org/archives/i18n/>.

To subscribe, send mail to the friendly person at `i18n-request at xfree86.org` explaining your motivation.

X11 fonts

Address: `fonts at xfree86.org`

This mailing list addresses the people who work on Unicode fonts and the font subsystem for the X11/XFree86 system.

Archives are at <http://devel.xfree86.org/archives/fonts/>.

The Unicode HOWTO

To subscribe, send mail to the overworked person at `fonts-request` at `xfree86.org` explaining your motivation.