

X<sub>Ǝ</sub>TeX:  
TeX *plus* Unicode *plus* OpenType...  
*minus* complexity

Jonathan Kew

SIL International

November 17, 2007

# History

**April 2004:** first public release of Xe<sub>Y</sub>TeX 0.3 (for Mac OS X only)

- built-in support for Unicode
- access to all fonts installed on the user's computer
- AAT (Apple Advanced Typography) for typographic features

# History

**April 2004:** first public release of Xe<sub>Λ</sub>TeX 0.3 (for Mac OS X only)

- built-in support for Unicode
- access to all fonts installed on the user's computer
- AAT (Apple Advanced Typography) for typographic features

**February 2005:** Xe<sub>Λ</sub>TeX reaches version 0.9

- OpenType layout features
- mature Λ<sub>Λ</sub>TeX support packages

# History

**April 2004:** first public release of Xe<sub>Y</sub>TeX 0.3 (for Mac OS X only)

- built-in support for Unicode
- access to all fonts installed on the user's computer
- AAT (Apple Advanced Typography) for typographic features

**February 2005:** Xe<sub>Y</sub>TeX reaches version 0.9

- OpenType layout features
- mature L<sup>A</sup>TeX support packages

**April 2006** (BachoTeX): Xe<sub>Y</sub>TeX for Linux (first public release)

**June 2006:** Xe<sub>Y</sub>TeX built for Windows by Akira Kakuto

# History

**April 2004:** first public release of Xe<sub>Λ</sub>TeX 0.3 (for Mac OS X only)

- built-in support for Unicode
- access to all fonts installed on the user's computer
- AAT (Apple Advanced Typography) for typographic features

**February 2005:** Xe<sub>Λ</sub>TeX reaches version 0.9

- OpenType layout features
- mature Λ<sub>Λ</sub>TeX support packages

**April 2006** (BachoTeX): Xe<sub>Λ</sub>TeX for Linux (first public release)

**June 2006:** Xe<sub>Λ</sub>TeX built for Windows by Akira Kakuto

**February 2007:** T<sub>Λ</sub>EX Live 2007 includes Xe<sub>Λ</sub>TeX for all platforms

- many thanks to Karl Berry, and all TL builders and testers!

## Character sets: **7 bits**, 8 bits, 16 bits, more...

The original T<sub>E</sub>X82 program used a 7-bit input text encoding, and supported 8-bit character codes internally for font access (up to 256 characters per font, though Knuth's own fonts only used 128).

A 7-bit input encoding, supplemented by conventions such as ``...'` for "...", `---` for —, etc., was designed primarily for English typesetting.

For European languages needing accented letters, etc., the Computer Modern fonts and Plain T<sub>E</sub>X supported various macros (`\'a` `\`e` `\u` `\i` `\^o` `\~u` etc.), which became a de facto standard for representing such letters.

Adaptations (e.g., Babel) for other languages sometimes added language-specific macros such as the active `"` character to provide more convenient input.

# Character sets: 7 bits, **8 bits**, 16 bits, more...

TeX 3.0 extended the program to have 8-bit character support throughout, including text input as well as font access.

Simplified the support of European languages, including input text using 8-bit national codepages, and hyphenation of words including accented characters.

For multilingual documents, even using Latin script, it can be difficult to find an 8-bit encoding that provides all the necessary characters.

Adding non-Latin scripts to the mix still leads to great complexity: a variety of custom encodings, fonts, pre-processors, macros, etc.

# Character sets: 7 bits, 8 bits, **16 bits**, more...

Unicode was initially conceived as a 16-bit character set, supporting up to 65,536 characters, that would include all the characters of all the world's scripts, eliminating the need for multiple codepages, and simplifying data interchange.

The major early effort to support Unicode in a T<sub>E</sub>X system was Omega, which was initially a 16-bit system (though it has since grown). However, Omega never really achieved stability and widespread use.

With time, it became clear that 65,536 characters are not sufficient, and Unicode was extended; it is now *incorrect* to think of Unicode as a 16-bit character set.



## Character sets: 7 bits, 8 bits, 16 bits, **more...**

The Unicode standard now defines around 100,000 distinct characters, and the number is still growing. The largest possible Unicode character code is U+10FFFF, which requires 21 bits to represent as a simple integer.

There are multiple ways to represent these codes in a file, with UTF-8 (a variable-length sequence of 8-bit bytes) usually being the most convenient.

One key feature of the UTF-8 encoding form for Unicode is that the 7-bit ASCII character set is unchanged; the exact same byte values are still used to represent these 128 characters.

Both Xe<sub>Y</sub>TeX and lua<sub>Y</sub>TeX now support the full Unicode character set. Although they use different internal representations, both read UTF-8 text files natively.

# Typesetting Unicode text with Xe<sub>Y</sub>TeX

Accented characters (many more than in any legacy codepage):

```
\font\txt="Charis SIL" at 14pt \txt
```

Óðinn átti tvá bræðr. Hét annarr Vé, en annarr Vílir.

Hej Slované, ještě naše slovanská řeč žije.

Dünyayı verelim çocuklara hiç değilse bir günlüğüne.

Kuř béga Šešùpè, kuř Nēmunas tēka, taĩ mūsų tėvynė, ...

# Typesetting Unicode text with X<sub>Y</sub>TeX

Accented characters (many more than in any legacy codepage):

```
\font\txt="Charis SIL" at 14pt \txt
Óðinn átti tvá bræðr. Hét annarr Vé, en annarr Vílir.
Hej Slované, ještě naše slovanská řeč žije.
Dünyayı verelim çocuklara hiç değilse bir günlüğüne.
Kuř béga Šešùpè, kuř Nēmunas tēka, taĩ mūsų tėvỹnè, ...
```

Óðinn átti tvá bræðr. Hét annarr Vé, en annarr Vílir.

Hej Slované, ještě naše slovanská řeč žije.

Dünyayı verelim çocuklara hiç değilse bir günlüğüne.

Kuř béga Šešùpè, kuř Nēmunas tēka, taĩ mūsų tėvỹnè, gražì  
Lietuvà.

# Complex scripts: Arabic, Indic, etc.

As well as the extended Latin character set, complex non-Latin scripts are supported through industry-standard Unicode encoding and fonts:

```
\fontspec{Charis SIL}Русский  
\fontspec{Devanagari MT}हिन्दी  
\fontspec{STKaiti}汉语  
\fontspec[Script=Arabic]{Noori Nastaliq MT}  
    \beginR نستعلیق اردو \endR  
\fontspec{AppleGothic}한국어  
\fontspec{Gentium}Ελληνικά  
\fontspec{Lucida Grande}עברית  
\fontspec{Hiragino Mincho Pro}日本語
```

# Complex scripts: Arabic, Indic, etc.

As well as the extended Latin character set, complex non-Latin scripts are supported through industry-standard Unicode encoding and fonts:

```
\fontspec{Charis SIL}Русский
```

*Русский*

```
\fontspec{Devanagari MT}हिन्दी
```

हिन्दी

```
\fontspec{STKaiti}汉语
```

汉语

```
\fontspec[Script=Arabic]{Noori Nastaliq MT}
```

اُردُو نستعلیق

```
\beginR نستعلیق اُردُو \endR
```

한국어

```
\fontspec{AppleGothic}한국어
```

Ελληνικά

```
\fontspec{Gentium}Ελληνικά
```

עברית

```
\fontspec{Lucida Grande}עברית
```

日本語

```
\fontspec{Hiragino Mincho Pro}日本語
```

# Far Eastern scripts (Chinese, Japanese, Korean)

They're just more characters; no special effort required:

# Far Eastern scripts (Chinese, Japanese, Korean)

They're just more characters; no special effort required:

```
\font\han="STSong" at 16pt
```

```
\font\rom="Gentium" at 8pt
```

```
\def\hc#1#2{\vbox{\hbox{\han #1}  
  \hbox{\kern10pt\rom #2}}}
```

```
\vbox{\hc{書<}{ka-ku}  
  \hc{最も}{motto-mo}  
  \hc{最後}{sai-go}  
  \hc{働<}{hatara-ku}  
  \hc{海}{umi}}
```

# Far Eastern scripts (Chinese, Japanese, Korean)

They're just more characters; no special effort required:

```
\font\han="STSong" at 16pt
```

```
\font\rom="Gentium" at 8pt
```

```
\def\hc#1#2{\vbox{\hbox{\han #1}  
  \hbox{\kern10pt\rom #2}}}
```

```
\vbox{\hc{書<}{ka-ku}  
  \hc{最も}{motto-mo}  
  \hc{最後}{sai-go}  
  \hc{働<}{hatara-ku}  
  \hc{海}{umi}}
```

書<  
ka-ku

最も  
motto-mo

最後  
sai-go

働<  
hatara-ku

海  
umi



# Vertical text with X<sub>Y</sub>T<sub>E</sub>X

```

\font\body="STKaiti:vertical" at 7pt \body
\font\bold="STHeiti:vertical" at 7pt
\font\ttitle="STHeiti:vertical" at 10pt
\centerline{\ttitle 三 国 演 义} \medskip
\centerline{\bold 〔明〕罗贯中} \smallskip
\leftline{词曰：}
滚滚长江东逝水，浪花淘尽英雄。是非成败转头空：青山依旧在，几
度夕阳红。
白发渔樵江渚上，惯看秋月春风。一壶浊酒喜相逢：古今多少事，都
付笑谈中。
\medskip
\centerline{\bold 第一回} \smallskip
\centerline{\bold 宴桃园豪杰三结义      斩黄巾英雄首立
功}\medskip
话说天下大势，分久必合，合久必分：周末七国分争，并入于秦；及
秦灭之后，
楚、汉分争，又并入于汉；汉朝自高祖斩白蛇而起义，一统天下，后
来光武中兴
，传至献帝，遂分为三国。推其致乱之由，殆始于桓、灵二帝。桓帝
禁锢善类，
崇信宦官。及桓帝崩，灵帝即位，大将军窦武、太傅陈蕃，共相辅
佐；时有宦官
曹节等弄权，窦武、陈蕃谋诛之，机事不密，反为所害，中涓自此愈
横。
\par

```

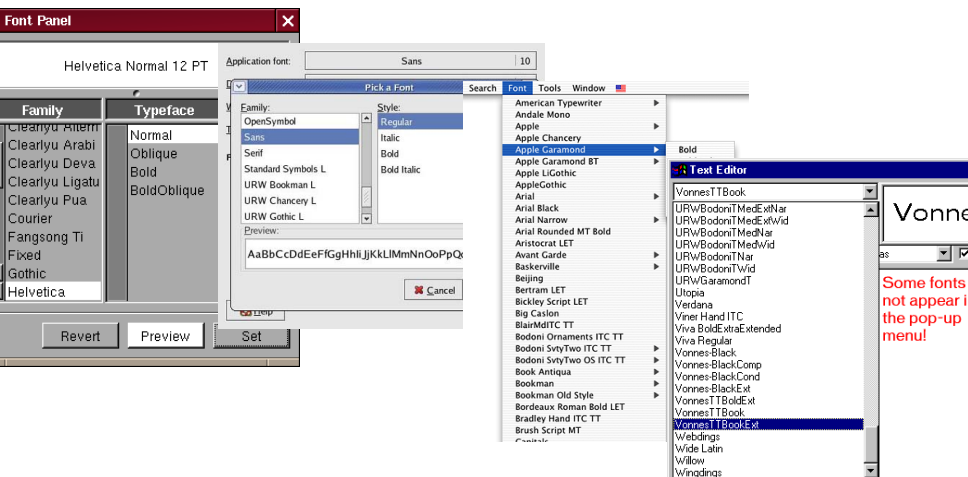


# Font access

How difficult should it be to specify the font you want to use?

# Font access

How difficult should it be to specify the font you want to use?



# Font access

How difficult should it be to specify the font you want to use?

```
\font\helv = phvr8t
```

```
\font\timesrom = ptmr8t
```

```
\font\timesital = ptmri8t
```

```
\font\charter = bchr8t at 12pt
```

# Font access

How difficult should it be to specify the font you want to use?

```
\font\helv = "Helvetica"
```

```
\font\timesrom = "Times New Roman"
```

```
\font\timesital = "Times New Roman Italic"
```

```
\font\charter = "Bitstream Charter" at 12pt
```

```
\font\garamond = "Adobe Garamond Pro"
```

```
\font\chancery = "Apple Chancery"
```

```
\font\dingbats = "Zapf Dingbats"
```

```
\font\japanese = "Hiragino Mincho Pro W3"
```

```
\font\chinese = "STKaiti" at 16pt
```

```
\font\arabic = "Scheherazade"
```

# Advanced features of modern fonts

Professional-quality OpenType fonts may include many features that were formerly available only as separate, custom-encoded ‘expert’ fonts, such as small capitals, lining and oldstyle numerals, extra ligatures and swash glyphs, etc.

# Advanced features of modern fonts

Professional-quality OpenType fonts may include many features that were formerly available only as separate, custom-encoded ‘expert’ fonts, such as small capitals, lining and oldstyle numerals, extra ligatures and swash glyphs, etc.

A font specification in Xe<sub>2</sub>TeX can include ‘tags’ that enable such features, if supported by the font:

```
\font\A="Adobe Garamond Pro"  
\A Adobe Garamond Pro: 12345  
\font\B="Adobe Garamond Pro:+smcp"  
\B Garamond with Small Caps  
\font\C="Adobe Garamond Pro:+onum"  
\C Oldstyle Numerals: 12345
```



# Advanced features of modern fonts

Professional-quality OpenType fonts may include many features that were formerly available only as separate, custom-encoded ‘expert’ fonts, such as small capitals, lining and oldstyle numerals, extra ligatures and swash glyphs, etc.

A font specification in Xe<sub>2</sub>TeX can include ‘tags’ that enable such features, if supported by the font:

```
\font\A="Adobe Garamond Pro"
```

```
\A Adobe Garamond Pro: 12345
```

```
\font\B="Adobe Garamond Pro:+smcp"
```

```
\B Garamond with Small Caps
```

```
\font\C="Adobe Garamond Pro:+onum"
```

```
\C Oldstyle Numerals: 12345
```

Adobe Garamond Pro: 12345

GARAMOND WITH SMALL CAPS

Oldstyle Numerals: 12345

# OpenType families with optical sizes

Where an OpenType font family includes optical size information, Xe<sub>Λ</sub>TeX will automatically load the appropriate face for the size being used:

```
\font\A="Briosio Pro" at 7pt  
\A This is \fontname\A
```

```
\font\B="Briosio Pro" at 11pt  
\B This is \fontname\B
```

```
\font\C="Briosio Pro" at 16pt  
\C This is \fontname\C
```

```
\font\D="Briosio Pro" at 24pt  
\D This is \fontname\D
```

# OpenType families with optical sizes

Where an OpenType font family includes optical size information, X<sub>Y</sub>T<sub>E</sub>X will automatically load the appropriate face for the size being used:

```
\font\A="Briosio Pro" at 7pt  
\A This is \fontname\A
```

```
\font\B="Briosio Pro" at 11pt  
\B This is \fontname\B
```

```
\font\C="Briosio Pro" at 16pt  
\C This is \fontname\C
```

```
\font\D="Briosio Pro" at 24pt  
\D This is \fontname\D
```

This is "Briosio Pro Caption" at 7.0pt

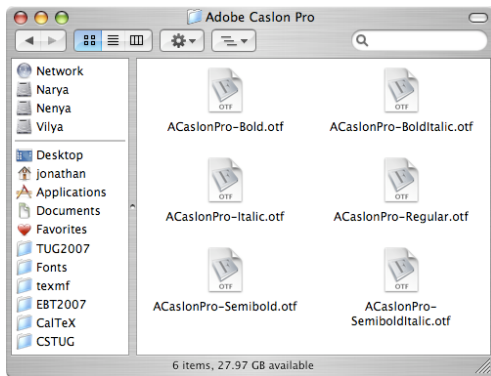
This is "Briosio Pro" at 11.0pt

This is "Briosio Pro  
Subhead" at 16.0pt

This is "Briosio  
Pro Display" at  
24.0pt

# Installing a new font

How about installing a new font, purchased or downloaded as an OpenType or TrueType file?



# Installing a new font

How about installing a new font, purchased or downloaded as an OpenType or TrueType file?

The traditional way:

- Rename font files according to a cryptic convention
- Possibly convert fonts to Type 1 (.pfb) format
- Generate metrics and virtual font descriptions
- Compile the metrics and virtual fonts into binary form
- Install fonts, metrics, virtual fonts into proper subdirectories of the texmf hierarchy
- Create map files for dvips, pdftex, dvipdfm, xdvi, ...
- Install map files and add entries to the drivers' configurations
- Update hash tables if installing system-wide

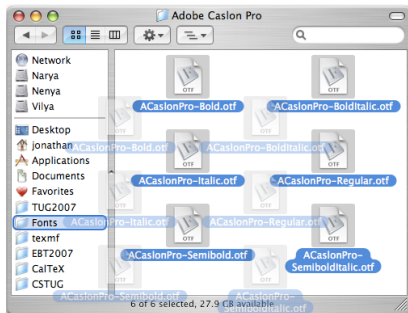
See the 100-page *Font Installation Guide* and other documents for more details.

# Installing a new font

How about installing a new font, purchased or downloaded as an OpenType or TrueType file?

The X<sub>Y</sub>TeX way:

- Drag the font files to the system's Fonts folder



- Ready to use via `\font\caslon = "Adobe Caslon Pro"`

# Using LaTeX with the XeTeX engine

Simply run `xelatex` instead of `pdflatex`.

# Using L<sup>A</sup>T<sub>E</sub>X with the Xe<sub>Y</sub>TeX engine

Simply run `xelatex` instead of `pdflatex`.

- normal L<sup>A</sup>T<sub>E</sub>X document classes and packages are used to format the document
- the key difference is that Unicode text and fonts can be used with no `\inputenc` package or special font setup



# Using Λ<sub>Λ</sub>TeX with the Xe<sub>Λ</sub>TeX engine

Simply run `xelatex` instead of `pdflatex`.

- normal Λ<sub>Λ</sub>TeX document classes and packages are used to format the document
- the key difference is that Unicode text and fonts can be used with no `\inputenc` package or special font setup

Λ<sub>Λ</sub>TeX is primarily a document markup and formatting language built on T<sub>Λ</sub>X; meanwhile, Xe<sub>Λ</sub>TeX extends the T<sub>Λ</sub>X engine to better handle Unicode and fonts. So they usually work well together.

# Using Λ<sub>Λ</sub>TeX with the Xe<sub>Λ</sub>TeX engine

Simply run `xelatex` instead of `pdflatex`.

- normal Λ<sub>Λ</sub>TeX document classes and packages are used to format the document
- the key difference is that Unicode text and fonts can be used with no `\inputenc` package or special font setup

Λ<sub>Λ</sub>TeX is primarily a document markup and formatting language built on T<sub>Λ</sub>TeX; meanwhile, Xe<sub>Λ</sub>TeX extends the T<sub>Λ</sub>TeX engine to better handle Unicode and fonts. So they usually work well together.

Common problem areas include:

- 8-bit text encodings supported via `inputenc` (convert to Unicode instead)
- custom-encoded fonts that are not Unicode-compliant
- packages that depend on a specific device driver (usually drawing packages that use PostScript output)

# Using L<sup>A</sup>T<sub>E</sub>X with the XeTeX engine

```
\documentclass{beamer}
\usetheme{Singapore}
\usefonttheme{serif} \usecolortheme{crane}

\usepackage{fontspec}
\setmainfont{Optima}
\setmonofont[Scale=MatchLowercase]{Andale Mono WT J}

\usepackage{xunicode,xltxtra,url,parskip,mdwlist}

\title{\XeTeX:\\
  $\phantom{\hbox{...}}$\TeX\ \emph{plus} Unicode
  \emph{plus} OpenType...\\ \emph{minus} complexity}
\author{Jonathan Kew}
\institute{SIL International}
\date{November 17, 2007}

\begin{document}
\maketitle
```



# fontspec: a simple interface to fonts

X<sub>Y</sub>T<sub>E</sub>X makes it easy to access a single font, at a specific size, by using the font name in a `\font` command.

However, L<sup>A</sup>T<sub>E</sub>X really needs more than this. To change the font of a L<sup>A</sup>T<sub>E</sub>X document, we may need to change not only the Roman but also the Bold and Italic faces, used for things like `\emph` or headings, at many different sizes (`\normalsize`, `\small`, `\footnotesize`, etc.)

Normally this is handled by a font-specific package such as `garamond.sty`, together with `.fd` files that link font family abbreviations, style options, encodings, and sizes to specific real or virtual fonts. Creating all this infrastructure accounts for much of the complexity of L<sup>A</sup>T<sub>E</sub>X font setup.

# fontspec: a simple interface to fonts

The fontspec package by Will Robertson eliminates all this complexity for Xe<sub>Y</sub>TeX users, by providing high-level font selection macros that rely on Xe<sub>Y</sub>TeX's capabilities behind the scenes.

# fontspec: a simple interface to fonts

The fontspec package by Will Robertson eliminates all this complexity for X<sub>Y</sub>TeX users, by providing high-level font selection macros that rely on X<sub>Y</sub>TeX's capabilities behind the scenes.

To change the typefaces used in a L<sup>A</sup>T<sub>E</sub>X document, a typical preamble can simply include a few lines such as:

```
\usepackage{fontspec}  
  
\setmainfont{Adobe Caslon Pro}  
\setsansfont{Lucida Sans}  
\setmonofont{Lucida Console}
```

(Usually, the option [Mapping=tex-text] should be added to the declarations, for compatibility with T<sub>E</sub>X typing conventions.)

# fontspec: a powerful interface to fonts

Many additional options are available to access advanced features of OpenType fonts. For example, it is possible to:

- specify default features to be applied to all fonts
- automatically scale typefaces to match ex-height, cap-height, etc.
- define new font family commands (like `\rmfamily`, `\sffamily`, etc.) that respect the L<sup>A</sup>T<sub>E</sub>X environment but access any available fonts
- precisely control which faces are used for Bold, Italic, etc., in a complex font family
- specify script- and language-specific options for fonts that support these features

See the extensive fontspec documentation for more details.



# xunicode for compatibility

To use LaTeX input conventions for accents and other ‘special’ characters with Unicode-encoded fonts, we need to redefine these macros to access the correct Unicode codepoints, otherwise LaTeX will try to print them from a variety of custom-encoded fonts.

Ross Moore’s package `xunicode` does this for many common LaTeX control sequences; it includes a couple of thousand definitions such as:

```
\DeclareEncodedCompositeCharacter{\UTFfencname}{\`}{0300}{02CB} % Combining grave accent
\DeclareEncodedCompositeCharacter{\UTFfencname}{\'}{0301}{02CA} % Combining acute accent
\DeclareEncodedCompositeCharacter{\UTFfencname}{\^}{0302}{02C6} % Combining circumflex accent
\DeclareEncodedCompositeCharacter{\UTFfencname}{\~}{0303}{02DC} % Combining tilde
\DeclareEncodedCompositeCharacter{\UTFfencname}{\=}{0304}{02C9} % Combining macron

\DeclareUTFcharacter[\UTFfencname]{x00A1}{\textexclamdown}
\DeclareUTFcharacter[\UTFfencname]{x00A2}{\textcent}
\DeclareUTFcharacter[\UTFfencname]{x00A3}{\textsterling}
\DeclareUTFcharacter[\UTFfencname]{x00A4}{\textcurrency}
\DeclareUTFcharacter[\UTFfencname]{x00A5}{\textyen}
```

# arabxetex

Package by François Charette providing an ArabTeX-like interface for typesetting languages in Arabic script with Xe<sub>Y</sub>L<sup>A</sup>T<sub>E</sub>X. Input text is either Unicode Arabic or a Latin transliteration.

# arabxetex

Package by François Charette providing an ArabTeX-like interface for typesetting languages in Arabic script with Xe<sub>La</sub>TeX. Input text is either Unicode Arabic or a Latin transliteration.

مِنَ الْقُرْآنِ الْكَرِيمِ، سُورَةُ السَّجْدَةِ ١٥-١٦:

إِنَّمَا يُؤْمِنُ بِآيَاتِنَا الَّذِينَ إِذَا ذُكِّرُوا بِهَا خَرُّوا سُجَّدًا وَسَبَّحُوا بِحَمْدِ رَبِّهِمْ وَهُمْ لَا يَسْتَكْبِرُونَ ﴿١٥﴾  
تَتَجَافَى جُنُوبُهُمْ عَنِ الْمَضَاجِعِ يَدْعُونَ رَبَّهُمْ خَوْفًا وَطَمَعًا وَمِمَّا رَزَقْنَاهُمْ يُنفِقُونَ ﴿١٦﴾

# arabxetex

Package by François Charette providing an ArabTeX-like interface for typesetting languages in Arabic script with X<sub>Y</sub>L<sup>A</sup>T<sub>E</sub>X. Input text is either Unicode Arabic or a Latin transliteration.

مِنْ الْقُرْآنِ الْكَرِيمِ، سُورَةُ السَّجْدَةِ ١٥-١٦:

إِنَّمَا يُؤْمِنُ بِآيَاتِنَا الَّذِينَ إِذَا ذُكِّرُوا بِهَا خَرُّوا سُجَّدًا وَسَبَّحُوا بِحَمْدِ رَبِّهِمْ وَهُمْ لَا يَسْتَكْبِرُونَ ﴿١٥﴾  
تَتَجَافَى جُنُوبُهُمْ عَنِ الْمَضَاجِعِ يَدْعُونَ رَبَّهُمْ خَوْفًا وَطَمَعًا وَمِمَّا رَزَقْنَاهُمْ يُنفِقُونَ ﴿١٦﴾

انبیاء بنی اسرائیل وچوں حضرت یونسؑ کجلیل القدر نبی ہن۔ آپ دا دُور نبوت ۸۲/۸۱ ق م کنوں ۵۳ ق م تک ہے۔ اوں ویلے مُملکت بنی اسرائیل دا بادشاہ یربعام ثانی ہا۔ آپ ایں مُملکت وچ ناصرت دے نال ”جات جفر“ دے مقام تے رَہندے ہن۔ اُج گل ایں جاہ کُوں ”خریۃ الزرع“ آہن۔ ایہ مقام جلیل جلیل دے مغربی پاسے ہارہاں میل اتین ناصرة دے شمال مشرقی پاسے ترائے میل دے فاصلے تے ہے۔ اٹھوں دے کھنڈرات دی کھدائی دے دوران پٹلگے جو ایہ مقام کانسی دے آخری دُور تقریباً ۱۵۵۰ ق م کنوں ۱۲۰۰ ق م تک آباد ہا۔ ولا دُوجھی دفعہ ایہ تقریباً ۱۲۰۰ ق م کنوں ۶۰۰ ق م تک آباد رہیا۔ حضرت یونسؑ ایں دُوجھے دُور وچ ہن۔ انہیں کھنڈرات کنوں ذرا شمالی پاسے کج وستی بھیندا ناں ہے ”مشهد“۔ روایت دے مطابق حضرت یونسؑ دا مزار اُتھائیں ہے۔

# jspacing, zhspacing

`jspacing`: package by Miyata Shigeru implementing detailed rules for Japanese character spacing (e.g., around punctuation marks).

ユニコードは、プラットフォームに係わらず、プログラムに係わらず、言語に係わらず、すべての文字に独立した番号を与えます。ユニコード標準は、アップル、ヒューレットパッカー、IBM、ジャストシステム、マイクロソフト、オラクル、SAP、サン、サイベースなどの産業界の主導的企業と他の多くの企業に採用されています。

# jspacing, zhspacing

`jspacing`: package by Miyata Shigeru implementing detailed rules for Japanese character spacing (e.g., around punctuation marks).

ユニコードは、プラットフォームに係わらず、プログラムに係わらず、言語に係わらず、すべての文字に独立した番号を与えます。ユニコード標準は、アップル、ヒューレットパッカード、IBM、ジャストシステム、マイクロソフト、オラクル、SAP、サン、サイベースなどの産業界の主導的企業と他の多くの企業に採用されています。

ユニコードは、プラットフォームに係わらず、プログラムに係わらず、言語に係わらず、すべての文字に独立した番号を与えます。ユニコード標準は、アップル、ヒューレットパッカード、IBM、ジャストシステム、マイクロソフト、オラクル、SAP、サン、サイベースなどの産業界の主導的企業と他の多くの企業に採用されています。

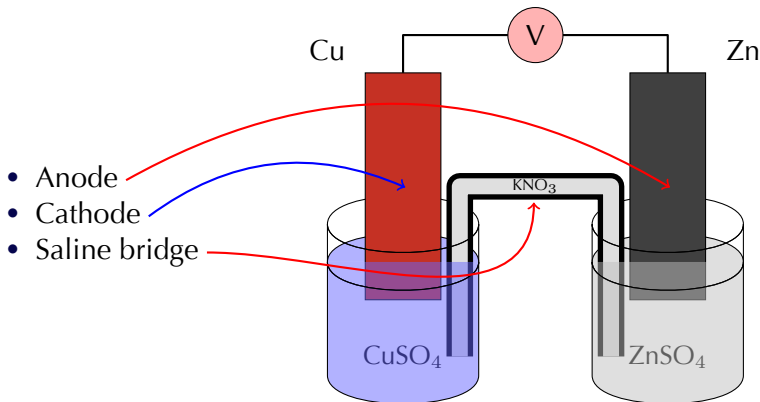
`zhspacing`: similar package for Chinese, by Yin Dian.

# Drawing packages such as *TikZ*

*TikZ* works because pgf supports the dvipdfm(x) driver:

# Drawing packages such as TikZ

TikZ works because pgf supports the dvipdfm(x) driver:



(example figure from <http://www.fauskes.net/nb/pgftikzexamples/>)



# driver-dependent packages

Many driver-dependent packages now support Xe<sub>Λ</sub>TeX, so that they work correctly regardless of the TeX engine used.

`graphics` supported via the `xetex.def` driver (by Ross Moore)

`color` supported via the `xetex.def` driver

`hyperref` recognizes Xe<sub>Λ</sub>TeX and uses appropriate `\specials`

`geometry` configuration file treats Xe<sub>Λ</sub>TeX like pdfTeX

`crop` configuration file treats Xe<sub>Λ</sub>TeX like pdfTeX

`pdfpages` supports the Xe<sub>Λ</sub>TeX primitives for PDF inclusion

`pgf/TikZ` configuration file treats Xe<sub>Λ</sub>TeX like dvipdfm

`pstricks` custom driver (by Miyata Shigeru) works with `xdvipdfmx`

# Using Xe<sub>Λ</sub>TeX with ConTeXt

texexec supports an `--xtx` option to use Xe<sub>Λ</sub>TeX in place of the default pdf<sub>Λ</sub>TeX engine. This allows the use of Xe<sub>Λ</sub>TeX font names and other options; see the [contextgarden.net](http://contextgarden.net) wiki.

# Using X<sub>Y</sub>T<sub>E</sub>X with ConT<sub>E</sub>Xt

texexec supports an `--xtx` option to use X<sub>Y</sub>T<sub>E</sub>X in place of the default pdfT<sub>E</sub>X engine. This allows the use of X<sub>Y</sub>T<sub>E</sub>X font names and other options; see the [contextgarden.net](http://contextgarden.net) wiki.

```
\definetypeface[Serapion][rm][Xserif][Serapion Pro]
\setupbodyfont[Serapion, 12pt]

\starttypescript[serif][didot][uc]
\definefontsynonym[DidotRegular]['Didot'] [encoding=uc]
\definefontsynonym[DidotItalic] ['Didot/I'] [encoding=uc]
\definefontsynonym[DidotBold] ['Didot/B'] [encoding=uc]
\definefontsynonym[DidotCaps] ['Didot;
    Letter Case=Small Capitals;Ligatures=!Common Ligatures']
    [encoding=uc]
\stoptypescript

\definedfont["Hoefer Text:mapping=tex-text;
    Style Options=Engraved Text;
    Letter Case=All Capitals;color=229966" at 24pt]
```

# Including PDF images

The primitive `\XeTeXpdffile` *[options]* "*filename*" includes a page from a PDF file.

*[option]* keywords include:

- page *number* (default is first page in the PDF)
- *crop*, *media*, *bleed*, *trim*, *art* select bounding box to use
- *scaled*, *xscaled*, *yscaled* *factor* (1000 = natural size)
- *width*, *height* *dimen*
- *rotated* *degrees*

From TeX's point of view, the image is like a big character; put it into an `\hbox` to measure height and width with macros.

# Including PDF images

The primitive `\XeTeXpdffile` *[options]* *"filename"* includes a page from a PDF file.

*[option]* keywords include:

- page *number* (default is first page in the PDF)
- *crop*, *media*, *bleed*, *trim*, *art* select bounding box to use
- *scaled*, *xscaled*, *yscaled* *factor* (1000 = natural size)
- *width*, *height* *dimen*
- *rotated* *degrees*

From TeX's point of view, the image is like a big character; put it into an `\hbox` to measure height and width with macros.

New in XeTeX 0.997: `\XeTeXpdfpagecount` *"filename"* accesses the number of pages in the given file; allows packages such as `pdfpages` to iterate over a complete file.

# JPEG, PNG, and BMP images

For raster images, use `\XeTeXpicfile` *[options]* *"filename"*

Options are the same as for PDF files, except for the PDF-specific bounding box and page selection options.

Again, XeTeX knows the dimensions of the image, so by putting it in a box and examining `\ht` and `\wd`, macros can do appropriate layout.

The L<sup>A</sup>T<sub>E</sub>X graphics package automatically recognizes XeTeX and uses these commands to handle pictures.

# JPEG, PNG, and BMP images

For raster images, use `\XeTeXpicfile [options] "filename"`

Options are the same as for PDF files, except for the PDF-specific bounding box and page selection options.

Again, XeTeX knows the dimensions of the image, so by putting it in a box and examining `\ht` and `\wd`, macros can do appropriate layout.

The LaTeX graphics package automatically recognizes XeTeX and uses these commands to handle pictures.

On Mac OS X, many additional image formats can be used with the `xdv2pdf` output driver, thanks to QuickTime, but for cross-platform compatibility JPEG and PNG are recommended.

## Support for EPS

It is also possible to use EPS images, because the `xdvipdfmx` driver can automatically convert these to PDF using Ghostscript or another distiller (see the `D` option in `dvipdfmx.cfg`).

Currently, this works with the L<sup>A</sup>T<sub>E</sub>X graphics package and the newest `xetex.def` driver definition. Other packages may need to be updated to accept EPS, as Xe<sub>Y</sub>TeX did not originally support this format.

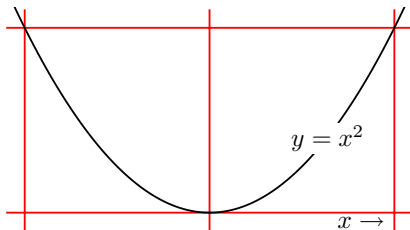


# Support for EPS

It is also possible to use EPS images, because the `xdvipdfmx` driver can automatically convert these to PDF using Ghostscript or another distiller (see the `D` option in `dvipdfmx.cfg`).

Currently, this works with the L<sup>A</sup>T<sub>E</sub>X graphics package and the newest `xetex.def` driver definition. Other packages may need to be updated to accept EPS, as Xe<sub>Y</sub>TeX did not originally support this format.

```
\includegraphics  
  [width=\hsize]  
    {combined.eps}
```



# TeX's math typesetting features

Essentially, a math formula is parsed into an *mlist* of various types of items: operators, variables, fractions, relations, etc. Then TeX converts this to nested hboxes and vboxes, using various font styles and positioning rules.

The basic operation of putting a character into a math formula is `\mathchar`, which takes a 16-bit value that is made up of a type (0–7), family (0–15), and character code (0–255). The family resolves to an appropriate font depending on `\textfont`, `\scriptfont`, and `\scriptscriptfont` settings.

For more complex elements, there are `\delimiter` and `\radical`, which take longer numbers with additional fields for small and large variants.

# TeX's math typesetting features

Math typesetting also depends on special font metrics in the ‘symbol’ and ‘extension’ fonts that include lists of variable-sized and extensible operators, delimiters, etc., in addition to the metrics of standard text fonts.

To allow characters to be entered directly in math formulas, instead of always using `\mathchar` or macros that expand to this, each character also has a `\mathcode` that determines what type of character it is, which math family it belongs to, etc. Similarly, there is `\delcode` for characters that are to behave as delimiters.

Math processing is described in detail in Appendix G of *The TeXbook*.

# Extending TeX for Unicode math

Unicode now provides standardized character codes for a huge number of math symbols, including styled letters (bold, fraktur, calligraphic, blackboard bold, etc.) that may be used for special purposes, as well as operators, relations, and so on.

The first step in allowing XeTeX to use Unicode for math typesetting is to extend the various math-related codes to support the full range of Unicode characters. To retain compatibility with old TeX macro files, this is done with new primitives `\XeTeXmathcode`, `\XeTeXmathchar`, `\XeTeXdelcode`, etc., that accept larger character codes.

At the same time, the number of permitted math font families has been increased from 16 to 256, to provide additional flexibility for macro package writers.

# Extending T<sub>E</sub>X for Unicode math

```
\XeTeXmathcode`\-= "2 "2 "2212 % minus sign
\xXeTeXmathcode`\ /= "0 "1 "2215 % division slash
\xXeTeXmathcode`\.= "0 "1 ` . % period
\xXeTeXmathcode`\,= "0 "1 ` , % comma
\xXeTeXmathcode`\Σ= "1 "2 `Σ % sum

\xXeTeXmathchardef\sum= "1 "2 `Σ
\xXeTeXmathchardef\prod= "1 "2 `Π
\xXeTeXmathchardef\intop= "1 "2 `∫
\xXeTeXmathchardef\ointop= "1 "2 `∮
\xXeTeXmathchardef\infty= "1 "2 `∞

\def\lceil{\XeTeXdelimiter"4 "3 "2308 } % ceiling
\def\rceil{\XeTeXdelimiter"5 "3 "2309 }
\def\lfloor{\XeTeXdelimiter"4 "3 "230A } % floor
\def\rfloor{\XeTeXdelimiter"5 "3 "230B }

\def\sqrt{\XeTeXradical"3 "221A }
```

# OpenType math fonts

With OpenType, it is possible for a single font to contain all the thousands of mathematical characters. In addition, a variety of metrics are needed to control the positioning and spacing of the parts of a formula.

Traditionally, T<sub>E</sub>X used a number of `\fontdimen` parameters from the `tfm` files for the math symbol and extension fonts. Equivalent information was simply not available with most other fonts.

With the release of Word 2007, Microsoft has defined a format for embedding such math-specific metrics in OpenType fonts. The first example of such a math-capable OpenType font is Cambria Math, included with a number of current MS products (including some freely-downloadable viewers).

# OpenType math fonts

Xe<sub>Y</sub>TeX is now able to use the OpenType ‘MATH’ table to control math typesetting, so that OpenType math fonts like Cambria can automatically be laid out correctly, with no additional metrics files.

$$f(x) = a_0 + \sum_{n=1}^{\infty} \left( a_n \cos \frac{n\pi x}{L} + b_n \sin \frac{n\pi x}{L} \right)$$

$$f(x) = a_0 + \sum_{n=1}^{\infty} \left( a_n \cos \frac{n\pi x}{L} + b_n \sin \frac{n\pi x}{L} \right)$$

# Handling text in native fonts

When using native system fonts (OpenType, AAT, Graphite), the minimal unit of text from the T<sub>E</sub>X point of view is the *native word* (a new node type).

Such a ‘native word’ might not correspond exactly to a natural-language word; it is simply a run of adjacent characters using a single font.

From T<sub>E</sub>X’s point of view, the word is like a box, with width, height, and depth, but its contents cannot be unboxed or otherwise inspected.

All layout within the word, like ligatures, kerning, or complex script behavior, is handled by a text layout engine (ATSUI, ICU Layout, or Graphite), which determines the actual glyphs to be rendered and their positioning.

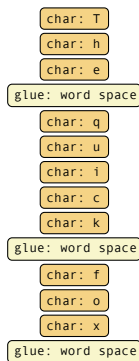


# Building a paragraph

TeX collects characters and spaces (glue) into a *horizontal list* which will be broken into lines.

# Building a paragraph

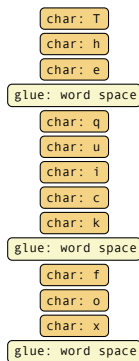
TeX collects characters and spaces (glue) into a *horizontal list* which will be broken into lines.



# Building a paragraph

TeX collects characters and spaces (glue) into a *horizontal list* which will be broken into lines.

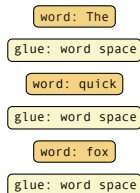
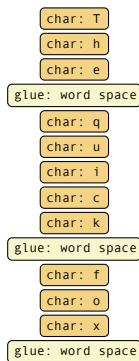
X<sub>La</sub>TeX gathers each run of characters in the same font into a *word*, so the horizontal list now consists of words and glue.



# Building a paragraph

TeX collects characters and spaces (glue) into a *horizontal list* which will be broken into lines.

XeTeX gathers each run of characters in the same font into a *word*, so the horizontal list now consists of words and glue.



# Hyphenation with system fonts

The horizontal list contains word nodes and glue:

Two glue different glue foxes

# Hyphenation with system fonts

The horizontal list contains word nodes and glue:

Two glue different glue foxes

If hyphenation is needed, the word node has to be deconstructed, and discretionary nodes inserted:

Two glue dif hyphen? fer hyphen? ent glue foxes

# Hyphenation with system fonts

The horizontal list contains word nodes and glue:

Two glue different glue foxes

If hyphenation is needed, the word node has to be deconstructed, and discretionary nodes inserted:

Two glue dif hyphen? fer hyphen? ent glue foxes

Now the line-break routine can use the discretionary breaks:

Two glue dif fer -  
ent glue foxes

Two differ-  
ent foxes

# Hyphenation with system fonts

The horizontal list contains word nodes and glue:

Two glue different glue foxes

If hyphenation is needed, the word node has to be deconstructed, and discretionary nodes inserted:

Two glue dif hyphen? fer hyphen? ent glue foxes

Now the line-break routine can use the discretionary breaks:

Two glue dif fer -  
ent glue foxes

Two differ-  
ent foxes

To avoid disrupting the low-level text layout, we must reassemble word fragments wherever a discretionary break was not used:

Two glue differ-  
ent glue foxes

Two differ-  
ent foxes



# Questions... and answers?

- Contact information
  - [mailto:jonathan\\_kew@sil.org](mailto:jonathan_kew@sil.org)
- Xe<sub>Y</sub>TeX web site and mailing list
  - <http://scripts.sil.org/xetex>
  - <http://tug.org/mailman/listinfo/xetex>
- Source repository for ongoing development
  - <http://scripts.sil.org/svn-public/xetex/TRUNK>

## Questions... and answers?

- Contact information
  - [mailto:jonathan\\_kew@sil.org](mailto:jonathan_kew@sil.org)
- X<sub>Y</sub>TeX web site and mailing list
  - <http://scripts.sil.org/xetex>
  - <http://tug.org/mailman/listinfo/xetex>
- Source repository for ongoing development
  - <http://scripts.sil.org/svn-public/xetex/TRUNK>

የኒኮድ ምንድን ነው? ما هي الشفرة الموحدة "يونيكود"؟ 什麼是Unicode(統一碼/標準萬國碼)? Što je Unicode? რა არის უნიკოდი? Τί είναι τὸ Unicode; מה זה יוניקוד? यूनिकोड क्या है? Hvað er Unicode? 유니코드とは何か? 유니코드에 대해? يونيكُد چیست؟ Что такое Unicode? Unicode ဖြစ်သလား? የኒኮድ ከንታይ ኢዩ?

