

## Typesetting Bangla in Unicode-aware T<sub>E</sub>X engines — from user experiences to development insights

Qutub Sajib

### Abstract

Typesetting Bangla (also known as Bengali) script in T<sub>E</sub>X was first introduced more than 15 years ago from now through a transliteration system. The system would require the user to input Bangla texts with fonts from Roman script in a specific ASCII transliteration scheme and then process the input to display the texts using its own METAFONT-generated Bangla fonts. The transliteration-based system falls short in terms of, among others, its harder-to-read source file and its requirement of one particular Bangla typeface family. With the introduction of Unicode-aware T<sub>E</sub>X engines, X<sub>Y</sub>T<sub>E</sub>X for example, and the emergence of Unicode-compliant free Bangla fonts, new possibilities have evolved. Today both X<sub>Y</sub>T<sub>E</sub>X and LuaT<sub>E</sub>X support Bangla typesetting allowing the user to input texts directly with Unicode Bangla fonts in the editor; thus making the source file easy-to-read and eliminating the need for transliteration schemes. Although several years have passed since the X<sub>Y</sub>T<sub>E</sub>X system was first seen to support Bangla typesetting, it is still in a state where the *finest* typographic quality is nearly unachievable for this particular script. While working with different Unicode Bangla fonts with both of these engines, several rendering issues were observed. In order to ensure fine typographic quality in Bangla script, these issues, among others, need attention from users and developers.

### 1 Introduction

The language Bangla, also known as Bengali, is the sixth most-spoken language in the world as reported by *Ethnologue* in its latest edition. Native speakers of this language are primarily from Bangladesh, a rather small country in south Asia. Another good number of native speakers are from West Bengal, a province of India. The Bangla script is one of the thirteen major Indic scripts and has made its way in the Unicode standard. Publishing in this language has a good history of many centuries. Like other Indic scripts, typesetting Bangla in T<sub>E</sub>X has seen several attempts in last few years but the typographic quality is still to find its prime.

Apart from the beautiful rendering of mathematical contents in T<sub>E</sub>X, another goal of this typesetting system was the *finest* typographic quality [1]. The same philosophy can be expected while typesetting other scripts including Bangla in T<sub>E</sub>X. Considering

the present-day situation of different attempts to typeset Bangla in T<sub>E</sub>X, this article tries to cover a few of them and provides some insights for future development.

### 2 Scope of this article

Before the Unicode standard was born for the purpose of writing most scripts of the world in computer, several attempts were taken to typeset Bangla in T<sub>E</sub>X using mostly ASCII based transliteration systems. Brief discussion of those systems are presented in Section 3. Following this, typesetting Bangla script in Unicode-aware T<sub>E</sub>X engines — X<sub>Y</sub>T<sub>E</sub>X and LuaT<sub>E</sub>X — are discussed with examples in subsequent sections. Some development ideas for this particular script are discussed in Section 11.

In this paper, when we say using X<sub>Y</sub>T<sub>E</sub>X we mean compiling the .tex file with the xelatex program; similarly when using LuaT<sub>E</sub>X we mean compiling the same file with lualatex. The T<sub>E</sub>X-specific examples presented in this paper were produced using the T<sub>E</sub>X Live 2019 distribution on a computer running GNU/Linux operating system, the Slackware 14.2 to be specific.

It is predictable that most Bangla documents would contain at least English, math and other scripts; but in this article we will consider typesetting only Bangla using two engines in T<sub>E</sub>X that supports the Unicode standard. This article does not cover the discussion about font selection techniques for different scripts except for Bangla. In order to learn selecting specific fonts for Roman (English) and maths along with Bangla, one may like to consult the fontspec package documentation.

### 3 An ASCII based transliteration system made the first attempt

As Anshuman Pandey mentioned [2], it was Avinash Chopde who first introduced Bengali (Bangla) in T<sub>E</sub>X in his itrans package. The subsequent developments in Bangla typesetting in T<sub>E</sub>X include the arosng by Muhammad Ali Masroor (later dropped), “ItxBengali” by Shrikrishna Patil, “Bengali Writer” and bwti by Abhijit Das, and the bengali along with a preprocessor beng by Pandey. Although the transliteration based system works, it has the limitation of its harder-to-read source file (Figure 1). Also, the font used to typeset the Bangla script is limited to one particular font family.

The bengali-omega package by Lakshmi K. Raut was probably the first to introduce typesetting Bangla using either the transliteration system or Unicode based direct input [4]. Although this package has the capability of taking Unicode Bangla as input,

```

কে লইবে মোর কার্য, কহে সন্ধ্যারবি
শুনিয়া জগৎ রহে নিরুত্তর ছবি।
মাটির প্রদীপ ছিল, সে কহিল, স্বামী
আমার যেটুকু সাধ্য করিব তা আমি।
-- রবীন্দ্রনাথ ঠাকুর

{\bn ke la\ibe mor kaarya, kahe sandhyaa rabi
"suni.yaa jagaT rahe niruttar chabi |
maa.tir pradiip chila, se kahila, sbaami
aamaar ye.tuku saadhya kariba taa ami |
-- rabindranaath .thaakur}

```

**Figure 1:** Typesetting Bangla in T<sub>E</sub>X with a transliteration based system as found in [2].

it uses the fonts as used in the package `bangtex` by Anshuman Pandey [3] to typeset the document.

#### 4 Then Unicode-aware T<sub>E</sub>X engines came into play

With the introduction of X<sub>Y</sub>T<sub>E</sub>X and LuaT<sub>E</sub>X engines, and the `fontspec` package for selecting TrueType and OpenType fonts, typesetting Bangla using Unicode fonts came into reality. Today, a good number of Unicode compliant Bangla fonts are freely available that seem to work with these engines.

To start, one needs a keyboard layout that would allow to input Unicode Bangla characters in the editor. In most Linux systems, a keyboard layout called `Probhat` is available to input Unicode Bangla characters. Another popular alternative is the *Avro Keyboard* (<https://www.omicronlab.com/index.html>) available for installation in MacOS X, Linux, and Windows operating systems. In `emacs`, three layouts are available, namely `bengali-inscript`, `bengali-ittrans`, and `bengali-probhat`. None of the Unicode Bangla keyboard layouts available today were designed with T<sub>E</sub>X users in mind; hence one may need to switch the layout frequently in order to type T<sub>E</sub>X-special characters (e.g., &, % etc.). Using any of these keyboard layouts, one can input Unicode compliant Bangla characters directly in the editor. An appropriate font containing the Bangla script has to be set via the `fontspec` package (details in Section 6). Then, upon processing the `.tex` file either with `xelatex` or `lualatex` one gets the typeset document.

Typesetting with the Unicode based system in T<sub>E</sub>X has the advantage over the older systems in terms of its easy-to-read source file. As shown in Figure 2, the source file can be read easily compared to the source file in Figure 1 for the same text. It should be mentioned here that to typeset the same text as in Figure 1, we have corrected some spellings following the original source [5]. But this system is not free from shortcomings; as seen in Figure 2 some codes are not readable which should read `"\vskip6pt\raggedleft"` (more on this later).

#### 5 Packages to support Unicode Bangla

The `polyglossia` package by François Charette is designed to provide support for typesetting Bangla

```

কে লইবে মোর কার্য, কহে সন্ধ্যারবি।
শুনিয়া জগৎ রহে নিরুত্তর ছবি।
মাটির প্রদীপ ছিল, সে কহিল, স্বামী,
আমার যেটুকু সাধ্য করিব তা আমি।
-- রবীন্দ্রনাথ ঠাকুর

\documentclass[utf8]{article}
\usepackage{fontspec}
\setmainfont{Noto Sans Bengali}
\textfont{Noto Sans Bengali}

```

**Figure 2:** Typesetting Bangla in X<sub>Y</sub>T<sub>E</sub>X with Unicode fonts.

script, along with many other scripts using Unicode fonts and Unicode-aware T<sub>E</sub>X engines. Although the package comes with language definitions files for different scripts including Bangla, it has many limitations. Even worse, the language definition file for Bangla script (`gloss-bengali.ldf`) that comes with this package has several words incorrectly spelled. It is well-known to the speakers and writers of Bangla language that the spelling of a few words has undergone changes and few are still in debate to agree on the “correct” spelling. But those misspelling in the file `gloss-bengali.ldf` are beyond any debate, as found in most dictionaries. For example, as of T<sub>E</sub>X Live 2019, the `gloss-bengali.ldf` file contains “সারপি”, “খঙ”, and few other words incorrectly spelled. (The author of this paper is not convinced to write incorrect spelling; so in order to see the *incorrectly* spelled words, one has to look into the above-mentioned file for the corresponding words).

The `polyglossia` package does come with few good features. It has the option to change the numbering style of L<sup>A</sup>T<sub>E</sub>X counters; i.e., one can typeset L<sup>A</sup>T<sub>E</sub>X counters with Bangla numerals. As the author of `latexbangla` stated, the `polyglossia` package’s current support does not offer any easy way to select fonts with Bangla script. This package tries to solve the problem by introducing control sequences to select several fonts for Bangla script but it has its shortcomings, too. As of our knowledge, there are no Unicode Bangla fonts designed to be used especially with T<sub>E</sub>X, the package sticks with the limited fonts available freely and tries to make bolded and monospaced texts using different fonts from different designers. This results into using the “AutoFakeBold” and “AutoFakeSlant” features of existing “Regular” style fonts. As a result, the typeset document becomes something passable but not of great aesthetic taste. The fonts used in `latexbangla` package are not available with T<sub>E</sub>X Live distribution. Besides, its dependence on the `ucharclasses` package makes it not usable with LuaT<sub>E</sub>X.

#### 6 Unicode compliant Bangla fonts to try with T<sub>E</sub>X

The T<sub>E</sub>X Live 2019 comes with the `gnu-freefont` package which contains Unicode fonts in both TTF and

দেশকালের বক্তৃতার জন্যই অভিকর্ষ—এ-কথা রবীন্দ্রনাথ এমন সময়ে বাঙালি-পাঠককে বলেছিলেন, যখন সে-সময়কে খুব বেশি ঔৎসুক্য এ-দেশে ছিল না। আসলে এ-দেশ হলো কবিতা আর গানের দেশ। রবীন্দ্রনাথ তাঁর অজস্র কবিতায় আর গানে হৃদ ও সুরের ঝংকারের প্রতি বাঙালি মনের চিরন্তন আকর্ষণকে রূপ দিয়েছেন; কিন্তু তাঁর নিজের মনটি যে আশ্চর্যজনকভাবে বিজ্ঞানানুগ ছিল—এ-খবর কখনে রাখেন?

দেশকালের বক্তৃতার জন্যই অভিকর্ষ—এ-কথা রবীন্দ্রনাথ এমন সময়ে বাঙালি-পাঠককে বলেছিলেন, যখন সে-সময়কে খুব বেশি ঔৎসুক্য এ-দেশে ছিল না। আসলে এ-দেশ হলো কবিতা আর গানের দেশ। রবীন্দ্রনাথ তাঁর অজস্র কবিতায় আর গানে হৃদ ও সুরের ঝংকারের প্রতি বাঙালি মনের চিরন্তন আকর্ষণকে রূপ দিয়েছেন; কিন্তু তাঁর নিজের মনটি যে আশ্চর্যজনকভাবে বিজ্ঞানানুগ ছিল—এ-খবর কখনে রাখেন?

**Figure 3:** Rendering of Bangla script in X<sub>Y</sub>T<sub>E</sub>X using Free Serif fonts: top: using MikT<sub>E</sub>X 2.8; bottom: using T<sub>E</sub>X Live 2019.

OTF formats and cover a wide range of Unicode character set. Fonts for Bangla script are available in serif and sans-serif version including regular and slanted styles. Unfortunately, no bold or bold italic fonts are available for Bangla script. Figure 3 shows rendering of few lines of Bangla script using Free Serif fonts with x<sub>Y</sub>l<sub>A</sub>TeX. In this figure, the typeset script on top is from the author’s own typesetting that was done sometime in 2009 using MikT<sub>E</sub>X 2.8 distribution. The same text when compiled with x<sub>Y</sub>l<sub>A</sub>TeX using T<sub>E</sub>X Live 2019, however, produces a different result. To be specific, few conjunct characters or ligatures are incorrectly rendered in T<sub>E</sub>X Live 2019. The reason behind this, whether with the fonts or the x<sub>Y</sub>l<sub>A</sub>TeX program, needs a fix in upcoming releases.

Besides the Free Serif fonts included in T<sub>E</sub>X Live 2019, Google noto fonts (<https://www.google.com/get/noto/>) include “Noto Serif Bengali” and “Noto Sans Bengali” in TTF format. The serif version comes with regular and bold styles while the sans-serif version contains other *seven* styles apart from the regular and bold. All these fonts can be used to typeset Unicode Bangla in T<sub>E</sub>X but they have to be downloaded and setup correctly so T<sub>E</sub>X finds them. The OTF version of Google noto fonts *are* available in T<sub>E</sub>X Live 2019 but they do not include the fonts for Bangla script.

There are a good number of Bangla fonts that can be freely downloaded from online. The Avro Keyboard website has a dedicated page to download fonts (<https://www.omicronlab.com/bangla-fonts.html>) and the Ekushey have their own page (<http://ekushey.org/index.php/page/33>). Most Unicode Bangla fonts available today, except the two fonts mentioned above, do not have their slanted/italic, bold etc. styles. Those fonts were not designed with professional publishing in mind, neither T<sub>E</sub>X users in mind. When using X<sub>Y</sub>T<sub>E</sub>X, users can use the “AutoFakeBold” and “AutoFakeSlant” features of fontspec package but the result is somewhat acceptable. In

LuaT<sub>E</sub>X, however, these features are not acceptable. So at this stage, users may have to stick with the limited styles available.

Besides the limited number of Unicode Bangla fonts or their availability in limited styles, rendering of Bangla characters in X<sub>Y</sub>T<sub>E</sub>X and LuaT<sub>E</sub>X needs deeper attention. In order to experiment with both of these engines, we selected three fonts to render the same texts. First one is the Free Serif fonts available in T<sub>E</sub>X Live 2019, second the “Noto Serif Bengali” fonts, and third the “Lohit Bengali” fonts. The last one of these three is available in most Linux distributions; in case it is not, can be downloaded from the Ekushey font page mentioned above. The font setup used in producing all examples in this article is given below. The “Noto Sans Bengali” is used in Figure 2 to show the verbatim texts. Considering the *x*-height of Free Serif fonts as “normal”, other fonts were scaled accordingly to get the identical output. This font setup works with both x<sub>Y</sub>l<sub>A</sub>TeX and lua<sub>Y</sub>l<sub>A</sub>TeX:

```
\usepackage{fontspec}
%
\newfontfamily{\freeserifbn}{FreeSerif.ttf}
[Script=Bengali, Ligatures=TeX]
\newfontfamily{\notoserifbn}
{NotoSerifBengali-Regular.ttf}
[Script=Bengali, Scale=0.85, Ligatures=TeX]
\newfontfamily{\notosansbn}
{NotoSansBengali-Light.ttf}
[Script=Bengali, Scale=0.85, Ligatures=TeX]
\ifxetex
\newfontfamily{\lohitbn}{lohit_bn.ttf}
[Path=/usr/share/fonts/TTF/,
Script=Bengali, Scale=0.82, Ligatures=TeX]
\else
\newfontfamily{\lohitbn}{lohit_bn.ttf}
[Script=Bengali, Scale=0.82, Ligatures=TeX]
\fi
```

## 7 The mystery of DOTTED CIRCLE

The DOTTED CIRCLE belongs to the “Geometric-Shapes” of the Unicode block and holds the character code U+25CC. It can be typeset with the T<sub>E</sub>X command `\char"25CC`, using the font setup mentioned above and compiling with either x<sub>Y</sub>l<sub>A</sub>TeX or lua<sub>Y</sub>l<sub>A</sub>TeX producing: ○ (here with Free Serif font). In many non-Roman scripts the DOTTED CIRCLE is used to indicate the position of various diacritical marks. In Bangla scripts, for example, short form of a vowel (known as *kaar*) replaces the actual vowel when the vowel comes after a consonant. Different *kaars* have their different positions to sit with consonants. The DOTTED CIRCLE is helpful to visualize the position of a *kaar* to sit with a consonant when the short form is typeset independently.



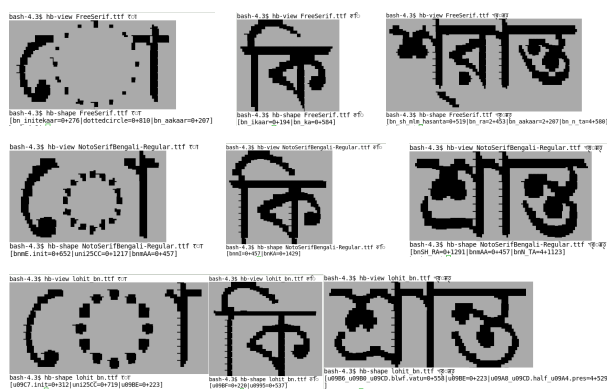
**Figure 4:** Typesetting short forms of vowels, independently or with consonants, produces different results in **xelatex** and **lualatex** for different fonts.

In Figure 4, all ten vowel short forms are shown independently (first row of a pair) and when they are combined with the consonant *kô* ক (second row of a pair). The example is shown in three pairs — using three fonts: Free Serif, Noto Serif Bengali, and Lohit Bengali — compiling with both **xelatex** and **lualatex**.

The problem arises when one wants to typeset short forms of vowels without the DOTTED CIRCLE. It is particularly necessary when children are taught these characters in schools. Ideally, in **T<sub>E</sub>X**-engines, one should be able to typeset the short forms of vowels. To our surprise, in its current situation, we cannot typeset them without the DOTTED CIRCLE with **xelatex**. We say *ideally* because of few reasons: (i) when a consonant takes a short form of vowel, the DOTTED CIRCLE disappears implying the presence of those characters independent of the DOTTED CIRCLE in current font; (ii) as described in the Unicode standard, the purpose of DOTTED CIRCLE is to show the relative position of the short form of a vowel; and (iii) when seen in a font viewer like **FontForge**, short forms of vowels can be found as independent glyphs. Therefore, rendering of vowel short forms in **xelatex** with a DOTTED CIRCLE can be considered as a bug.

## 8 Bangla script in HarfBuzz and Word processors

To understand the mystery of DOTTED CIRCLE in **T<sub>E</sub>X**, we can have a look in the output produced by the text rendering stack **HarfBuzz** as this program works behind the **X<sub>T<sub>E</sub>X</sub>** engine as well as in many Word processors. Two modules **hb-view** and **hb-shape** are available as part of **HarfBuzz** Indic Shaper and can be used in Linux shell to get the rendered output of a Unicode script. As shown in Figure 5, **HarfBuzz** produces the same result for different fonts as we



**Figure 5:** Rendering of Bangla script in **HarfBuzz** using different fonts in Linux shell.



**Figure 6:** Rendering of Unicode Bangla script in Word processors.

have seen with **xelatex**. The same result is also found in LibreOffice, OpenOffice, and MS Word (Figure 6).

From these examples it can be said that the **HarfBuzz** program is responsible for the unusual rendering of DOTTED CIRCLE in **X<sub>T<sub>E</sub>X</sub>** when typeset with the short forms of vowels. This assumption is also supported by the fact that the **LuaT<sub>E</sub>X** engine does not depend on **HarfBuzz** to render Unicode supported scripts and it produces expected results in terms of DOTTED CIRCLE rendering: (first row of second pair in Figure 4).

## 9 Dealing with hyphenation

In modern-day Bangla publications hyphenation is hardly seen, either because of technical limitations or lack of interest. Fortunately, **polyglossia** package supports hyphenation for Bangla script which works with both **xelatex** and **lualatex**. But unfortunately, the hyphenation rules in this package are too simple for this script. For example, one may like the word “একখানি” to be hyphenated between এক and খানি instead of what we see in Figure 7 (right).

## 10 Working with colors

Use of colors in texts can significantly improve the visual as well as readability for particular types of contents. Colorful texts can be essential in books written for children. For Bangla script, sometimes it is desirable to typeset different parts of a conjunct character in different colors to help children learn them separately. A good example can be to flag a

আমাদের পোস্টমাস্টার কলিকাতার হেলে। জলের মাহকে ডাঙায় তুলিলে বেরকম অবস্থা হয় এই গুণ গ্রামের মধ্যে আসিয়া পোস্টমাস্টারেরও সেই দশা উপস্থিত হইয়াছে। একখানি অঙ্ককার আটচালার মধ্যে তাঁকহার অফিস। অদূরে এটি পানাপুকুর একবং তাহার চার পাড়ে জঙ্গল।

**Figure 7:** Hyphenation supported by the `polyglossia` package works with both `xelatex` and `lualatex`.

কাকা যায়। ডাব খায়।	কাকা যায়।	xelatex
মৌরি রাখি কৌটা ভরি।	মৌরাখি	xelatex
	মৌরাখি	lualatex
		xelatex
		lualatex

**Figure 8:** Using different colors for different parts of a conjunct or ligature is a challenge; left: example from online (<http://tiny.cc/4xly9y>) right:  $\text{\TeX}$ -output.

*kaar* in a different color when it sits with a consonant, as can be seen in textbooks for children (Figure 8, top left). The same result in  $\text{\TeX}$ , either with  $\text{\XeTeX}$  or  $\text{\LuaTeX}$  is currently not achievable. The same problem arrives when trying to typeset a ligature flagging its different parts with different colors.

## 11 What is next?

To address the *finest* typographic quality in Bangla script, several things can be taken into consideration. A font family can be designed with  $\text{\TeX}$  users in mind and a supporting macro package can be developed. Use of only the `fontspec` package at the primary stage would be a good idea; integration with the `polyglossia` package may come next.

To solve the rendering issues especially with the DOTTED CIRCLE in  $\text{\XeTeX}$  and ligatures in  $\text{\LuaTeX}$ , attention could be made with the `HarfBuzz` developers. For now, because there are limitations in both  $\text{\XeTeX}$  and  $\text{\LuaTeX}$ , it is probably good to experiment with both of these engines. Eventually we may want to set up with one particular engine. This could take us to the  $\text{\LuaTeX}$  as this engine is supposed to be the successor of  $\text{\pdfTeX}$ .

Until a new dedicated font family is designed for the Bangla script, the Noto Bengali (both serif and sans-serif) fonts can be included in the future versions of  $\text{\TeX}$  Live. This would allow the users to try Unicode-aware  $\text{\TeX}$  engines with at least two font families including the already existing Free Serif and Free Sans fonts.

A keyboard layout could be designed for the purpose of Unicode Bangla character input making it `emacs`- and  $\text{\TeX}$ -friendly. In designing such a layout, the  $\text{\TeX}$  escape character (“\”), percent (“%”), ampersand (“&”) etc. keys can be left as is, so that these keys can be used to format Bangla texts without having to switch the keyboard layout.

## 12 Conclusion

The present-day Bangla publishing industry is mostly not using the fine typographic power of  $\text{\TeX}$ . The reasons behind this are many of which few are discussed here. Interests in solving those issues have been seen in recent years. Although usage of  $\text{\TeX}$  in Bangla fiction books could be a bigger challenge mostly due to non- $\text{\TeX}$  reasons, a good number of science books could be expected being typeset in  $\text{\TeX}$ . For this to happen, the current limitations of Unicode engines and fonts need to be addressed. The few insights we were able to uncover in this article could lead us to the beginning of the *finest* typographic quality in Bangla publishing.

## References

- [1] D. E. Knuth. *The  $\text{\TeX}$ book*, vol. A of *Computers & Typesetting*. Addison–Wesley Publishing Company, Massachusetts, 2012.
- [2] A. Pandey. Typesetting Bengali in  $\text{\TeX}$ . *TUGboat* 20(2):119–126, 1999.
- [3] A. Pandey. *Bengali for  $\text{\TeX}$* , 2.0 edition, 2002.
- [4] L. K. Raut. *Typesetting Bengali in  $\Omega$  using Velthuis Transliteration or Unicode Text*, 2006.
- [5] R. Tagore. The complete works of Rabindranath Tagore. Retrieved: 15 July 2019. <http://tiny.cc/tagore>

◇ Qutub Sajib  
China University of Geosciences,  
Wuhan  
388 Lumo Road, Wuhan 430074  
China  
[qsajib71@gmail.com](mailto:qsajib71@gmail.com)