

Experimental unicode mathematical typesetting: The unimath package

Will Robertson

2006/02/20 v0.01

Contents

1	Introduction	1			
2	Current NFSS methods	1			
3	Specification	2			
3.1	Dealing with real life	3	5.1	Enlarging the number of maths families	5
			5.2	<code>\DeclareMathSymbol</code> for unicode ranges	6
			5.3	User interface to <code>\DeclareSymbolFont</code>	7
			5.4	Setting up the ASCII ranges	8
I	The <code>unimath</code> package	3	6	Symbol definitions	9
4	Trying to understand \LaTeX	4			
5	This package	5	II	STIX table data extraction	9

1 Introduction

This document describes the `unimath` package, which is an *experimental* implementation of a macro to unicode glyph encoding for mathematical characters. Its intended use is for \XeTeX , although it is conjectured that small effect needs to be spent to create a cross-format package that would also work with \TeX .

As of \XeTeX v. 0.995, maths characters can be accessed in unicode ranges. Now, a proper method must be invented for real unicode maths support. Before any code is written, I'm writing a specification in order to work out what is required. Fairly significant pieces of the NFSS may have to be re-written, and I'm a little unsure where to start.

2 Current NFSS methods

In the following, $\langle NFSS\ decl. \rangle$ stands for something like `{T1}{lmr}{m}{n}`.

Maths symbol fonts Fonts for symbols: α , \leq , \rightarrow

`\DeclareSymbolFont{<name>}{NFSS decl.}`

Declares a named maths font such as `operators` from which symbols are defined with `\DeclareMathSymbol`.

Maths alphabet fonts Fonts for $ABC-xyz$, $\mathfrak{ABC}-\mathcal{XYZ}$, etc.

`\DeclareMathAlphabet{<cmd>}{NFSS decl.}`

For commands such as `\mathbf`, accessed through maths mode that are unaffected by the current text font, and which are used for alphabetic symbols in the ASCII range.

`\DeclareSymbolFontAlphabet{<cmd>}{<name>}`

Alternative (and optimisation) for `\DeclareMathAlphabet` if a single font is being used for both alphabetic characters (as above) and symbols.

Maths ‘versions’ Different maths weights can be defined with the following, switched in text with the `\mathversion{<maths version>}` command.

`\SetSymbolFont{<name>}{<maths version>}{NFSS decl.}`

`\SetMathAlphabet{<cmd>}{<maths version>}{NFSS decl.}`

Maths symbols Symbol definitions in maths for both characters (=) and macros (`\eqdef`): `\DeclareMathSymbol{<symbol>}{<type>}{<named font>}{<slot>}`
This is the macro that actually defines which font each symbol comes from and how they behave.

Delimiters, accents, and radicals are not dealt with yet.

3 Specification

In the ideal case, a single unicode font will contain all maths glyphs we need. Barbara Beeton’s STIX table provides the mapping between unicode maths glyphs and macro names (all 3298 of them!). A single command

`\setmathsfont[]{}`

would implement this for every every symbol and alphabetic variant. That means `\alpha` to α , `\leq` to \leq , etc., `\mathcal{H}` to \mathcal{H} and so on, all for unicode glyphs within a single font.

Furthermore, this package should deal well with unicode characters for maths input, as well. This includes using literal Greek letters in formulae, resolving to upright or italic depending on preference. This, and alphabetic variants via such commands as `\mathcal`, will be dealt with via X_YTeX’s ‘last minute’ font mapping features. (Or maybe not!)

Finally, maths versions must also be provided for. While I guess version selection in L^AT_EX will remain the same, the specification for choosing the version fonts will probably be an optional argument:

`\setmathsfont[version=bold,]{}`

All instances of ‘maths’ in command names will be aliased to ‘math’ for our American (or abbreviatory-minded) friends. Instances above of

`[]{}`

follow from my `fontspec` package, and therefore any additional `` specific to maths fonts will hook into `fontspec`’s methods.

3.1 Dealing with real life

Let's face it; there will probably be few cases where a single unicode maths font suffices. The upcoming STIX font comes to mind as a notable exception. It will therefore be necessary to delegate specific unicode ranges of glyphs to separate fonts.

At the lowest level, it will probably be necessary on occasion to simply use just one or two glyphs from another font; either because they look better or they're simply unavailable in the default font in use. This doesn't really require anything that won't already exist; a command analogous to `\DeclareMathSymbol` that accepts unicode *⟨slot⟩* ranges.

More generally, it would be nice to be able to say

`\setmathsfont [range=⟨unicode range⟩,⟨font features⟩] {⟨font name⟩}`

where *⟨unicode range⟩* is a comma-separated list of unicode slots and ranges such as {27D0–27EB, 27FF, 295B–297F}. Furthermore, preset names ranges could be used, such as `MiscMathSymbolsA`, with such ranges based on unicode chunks. The amount of optimisation required here to achieve acceptable performance has yet to be determined. Techniques such as saving out unicode subsets based on *⟨unicode range⟩* data to be `\input` in the next \LaTeX run are certainly a possibility.

File I

The **unimath** package

This is the package.

```
1 \ProvidesPackage{unimath}
2   [2006/02/20 v0.01 Unicode maths definitions]

   Things we need:
3 \newcounter{um@fam}
4
5 %% Kees can der Laan's simplification of Van der Groot's loop:
6 \def\um@Loop#1\um@Pool{#1\um@Loop#1\um@Pool}
7 \def\um@Break#1\um@Pool{}
8
9 \long\def\um@FOR #1 = [#2:#3] #4{%
10  {\csname#1\endcsname =#2\relax
11   \um@Loop #4%
12   \expandafter\advance\csname#1\endcsname\@ne
13   \expandafter\ifnum\csname#1\endcsname>#3\relax
14   \expandafter\um@Break
15   \fi
16   \um@Pool}}
17 \RequirePackage{fontspec}
```

Test: $\infty \angle \mathbb{Z} \mathfrak{A} \mathbb{H} \mathbb{I} \mathbb{J} \wedge \vee \cap \cup \int \iint \iiint$

4 Trying to understand L^AT_EX

Here's L^AT_EX's definition of `\DeclareMathSymbol`. Let's try and make sense of it.

```
\DeclareMathSymbol      #1: Symbol, e.g., \alpha or 'a'
                        #2: Type, e.g., \mathalpha
                        #3: Math font name, e.g., operators
                        #4: Slot, e.g., F1
                        <*>neveroutput>

18 \def\DeclareMathSymbol#1#2#3#4{%
```

First ensure the math font (e.g., operators) exists:

```
19   \expandafter\in@\csname sym#3\expandafter\endcsname
20   \expandafter{\group@list}%
21   \ifin@
```

Convert the slot number to two hex digits stored in `\count\z@` and `\count\tw@`, respectively:

```
22   \begingroup
23   \count\z@=#4\relax
24   \count\tw@\count\z@
25   \divide\count\z@\sixt@@n
26   \count@\count\z@
27   \multiply\count@\sixt@@n
28   \advance\count\tw@-\count@
```

The symbol to be defined can be either a command (`\alpha`) or a character (`a`). Branch for the former:

```
29   \if\relax\noexpand#1% is command?
30   \edef\reserved@a{\noexpand\in@{\string\mathchar}{\meaning#1}}%
31   \reserved@a
```

If the symbol command definition contains `\mathchar`, then we can provide the info that a previous symbol definition is being overwritten:

```
32   \ifin@
33   \expandafter\set@mathsymbol
34   \csname sym#3\endcsname#1#2%
35   {\hexnumber@\count\z@}\hexnumber@\count\tw@}%
36   \@font@info{Redeclaring math symbol \string#1}%
```

Otherwise, throw an error if the command name is already taken by a non-symbol definition:

```
37   \else
38   \expandafter\ifx
39   \csname\expandafter\@gobble\string#1\endcsname
40   \relax
41   \expandafter\set@mathsymbol
42   \csname sym#3\endcsname#1#2%
43   {\hexnumber@\count\z@}\hexnumber@\count\tw@}%
44   \else
45   \@latex@error{Command '\string#1' already defined}\@eha
46   \fi
47   \fi
```

And if the symbol input is a character:

```
48     \else
49     \expandafter\set@mathchar
50     \csname sym#3\endcsname#1#2
51     {\hexnumber@{\count\z@}\hexnumber@{\count\tw@}}}%
52     \fi
53 \endgroup
```

Everything previous was skipped if the maths font doesn't exist in the first place:

```
54 \else
55 \latex@error{Symbol font `#3' is not defined}\@eha
56 \fi}
```

The final macros that actually define the maths symbol with T_EX primitives. If the symbol definition is for a macro:

```
57 \def\set@mathsymbol#1#2#3#4{%
58 \global\mathchardef#2"\mathchar@type#3\hexnumber@#1#4\relax}
```

Or if it's for a character:

```
59 \def\set@mathchar#1#2#3#4{%
60 \global\mathcode`#2="\mathchar@type#3\hexnumber@#1#4\relax}
</neveroutput>
```

Summary For symbols, something like:

```
\def\DeclareMathSymbol#1#2#3#4{%
\global\mathchardef#1"\mathchar@type#2
\expandafter\hexnumber@\csname sym#2\endcsname
{\hexnumber@{\count\z@}\hexnumber@{\count\tw@}}}
```

For characters, something like:

```
\def\DeclareMathSymbol#1#2#3#4{%
\global\mathcode`#1"\mathchar@type#2
\expandafter\hexnumber@\csname sym#2\endcsname
{\hexnumber@{\count\z@}\hexnumber@{\count\tw@}}}
```

5 This package

We need to both redefine `\DeclareMathSymbol` to deal with unicode slots, as well as `\DeclareSymbolFont` to deal with 8-bit family numbers.

5.1 Enlarging the number of maths families

To start with, we've got a power of two as many `\fams` as before. So (from `lTFSSbas.dtx`) we want to redefine

```
61 \def\new@mathgroup{\alloc@8\mathgroup\chardef\@ccclvi}
62 \let\newfam\new@mathgroup
```

This is sufficient for L^AT_EX's `\DeclareSymbolFont,...`, commands to be able to define 256 named maths fonts. All we need now is a new `\DeclareMathSymbol`.

```

63 (*□□□□□□□)
64 \newfam\mta\newfam\mtb\newfam\mtc\newfam\mtd
65 \newfam\mte\newfam\mtf\newfam\mtg\newfam\mth
66 \newfam\mti\newfam\mtj\newfam\mtk\newfam\mtl
67 \newfam\mtm\newfam\mtn\newfam\mto\newfam\mtp
68 \newfam\mtq\newfam\mtr\newfam\mts\newfam\mtt
69 \newfam\mtu
70 (/□□□□□□□)

```

`\mtu`: math fam 24 of 255.

5.2 `\DeclareMathSymbol` for unicode ranges

This is mostly an adaptation from L^AT_EX's definition.

```

71 \def\DeclareUnicodeMathSymbol#1#2#3#4{%

```

First ensure the math font (e.g., operators) exists:

```

72 \expandafter\in@\csname sym#3\expandafter\endcsname
73 \expandafter{\group@list}%
74 \ifin@

```

No longer need to perform the obfuscated hex conversion, since X_ET_EX_{mathchar} (and friends) has a more simplified input than T_EX's `\mathchar`.

```

75 \begingroup

```

The symbol to be defined can be either a command (`\alpha`) or a character (`a`). Branch for the former:

```

76 \if\relax\noexpand#1% is command?
77 \edef\reserved@a{\noexpand\in@{\string\XeTeXmathchar}{\meaning#1}}%
78 \reserved@a

```

If the symbol command definition contains `\XeTeXmathchar`, then we can provide the info that a previous symbol definition is being overwritten:

```

79 \ifin@
80 \expandafter\set@xmathsymbol
81 \csname sym#3\endcsname#1#2{#4}%
82 \@font@info{Redeclaring math symbol \string#1}%

```

Otherwise, overwrite it if the symbol command definition contains plain old `\mathchar`:

```

83 \else
84 %\edef\reserved@a{\noexpand\in@{\string\mathchar}{\meaning#1}}%
85 %\reserved@a
86 %\ifin@
87 % \expandafter\set@xmathsymbol
88 % \csname sym#3\endcsname#1#2{#4}%

```

Otherwise, throw an error if the command name is already taken by a non-symbol definition:

```

89 %\else
90 %\expandafter\ifx

```

```

91         %\csname\expandafter\@gobble\string#1\endcsname
92         %\relax
93         \expandafter\set@xmathsymbol
94         \csname sym#3\endcsname#1#2{#4}%
95     %\else
96     % \latex@error{Command '\string#1' already defined}\@eha
97     %\fi
98     %\fi
99     \fi

```

And if the symbol input is a character:

```

100     \else
101     \expandafter\set@xmathchar
102     \csname sym#3\endcsname#1#2{#4}%
103     \fi
104     \endgroup

```

Everything previous was skipped if the maths font doesn't exist in the first place:

```

105     \else
106     \latex@error{Symbol font `#3' is not defined}\@eha
107     \fi}

```

The final macros that actually define the maths symbol with X_YTEX primitives. If the symbol definition is for a macro:

```

108 \def\set@xmathsymbol#1#2#3#4{%
109   \global\XeTeXextmathchardef#2"\mathchar@type#3 #1 #4\relax}

```

Or if it's for a character:

```

110 \def\set@xmathchar#1#2#3#4{%
111   \global\XeTeXextmathcode`#2="\mathchar@type#3 #1 #4\relax}

```

[For later] or if it's for a character code:

```

112 \def\DeclareUnicodeMathCode#1#2#3#4{%
113   \expandafter\set@xmathcode
114   \csname sym#3\endcsname{#1}{#2}{#4}}
115 \def\set@xmathcode#1#2#3#4{%
116   \global\XeTeXextmathcode#2="\mathchar@type#3 #1 #4\relax}
117 (*□□□□□□□)
118 \zf@fontspec{}{Cambria Math}
119 \DeclareSymbolFont{test}{EU1}{CambriaMath(0)}{m}{n}
120 \DeclareUnicodeMathSymbol{\infinity}{\mathord}{test}{"221E}
121 \DeclareUnicodeMathCode{65}{\mathalpha}{test}{119860}
122 (/□□□□□□□)

```

Test infinity: A^∞

5.3 User interface to \DeclareSymbolFont

\setmathfont [#5]: font features

#6: font name

Use fontspec to select a font to use.

```

123 \newcommand\setmathfont [2] [] {%
124   \zf@fontspec{#1}{#2}%

```

We need to hook into `fontspec` here to check if a family is loaded twice. This might be important if loading lots of individual glyphs.

```
125 \stepcounter{um@fam}%
126 \DeclareSymbolFont{um@fam\theum@fam}{EU1}{\zf@family}{\mddefault}{\updefault}%
```

Now when the list of unicode symbols is input, we want a suitable definition of its internal macro.

```
127 \def\unicode@math@symbol##1##2##3##4{%
128 \DeclareUnicodeMathSymbol{##2}{##3}{um@fam\theum@fam}{##1}}%
```

And now we input every single maths char. See File II for the source to `stix-tex.tex`.

```
129 \um@text@input{um@fam\theum@fam}%
130 \input stix-tex.tex}

131 (*□□□□□□□)
132 \setmathfont{Cambria Math}
133 \DeclareUnicodeMathSymbol{\testalef}{\mathord}{um@fam\theum@fam}{2135}
134 (/□□□□□□□)
```

Test aleph: \aleph

$\int Ax \stackrel{\text{def}}{=} \nabla \times Z$

5.4 Setting up the ascii ranges

We want it to be convenient for users to actually type in maths. The ASCII Latin characters should be used for italic maths, and the text Greek characters should be used for upright/italic (depending on preference) Greek, if desired.

```
\um@mathcode@offset This is a wrapper macro to save space:
135 \newcommand\um@mathcode@offset[3]{%
136 \DeclareUnicodeMathCode
137 {\numexpr\the\@tempcnta+#2\relax}
138 {\mathalpha}{#1}
139 {\numexpr\the\@tempcnta+#3\relax}%
140 }
```

`\um@text@input` And here're the text input to maths output mappings, wrapped up in a macro.

```
141 \newcommand\um@text@input[1]{%
Latin alphabet, uppercase and lowercase respectively:
142 \um@FOR @tempcnta = [0:25] {%
143 \um@mathcode@offset{#1}{65}{119860}%
144 \um@mathcode@offset{#1}{97}{119886}%
145 }%
```

Filling a hole for 'h', which maps to U+210E: PLANCK CONSTANT instead of the expected U+1D455: MATHEMATICAL ITALIC SMALL H (which is not assigned):

```
146 \DeclareUnicodeMathCode
147 {104}{\mathalpha}{#1}{8462}%
```


Greek alphabet, uppercase (note the hole after U+03A1: GREEK CAPITAL LETTER RHO):

```

148 \um@FOR @tempcnta = [0:23] {%
149   \DeclareUnicodeMathCode
150     {\ifnum \@tempcnta>16
151       \numexpr\the\@tempcnta+913\relax
152     \else
153       \numexpr\the\@tempcnta+913+1\relax
154     \fi}
155   {\mathalpha}{#1}
156   {\numexpr\the\@tempcnta+120546\relax}%

```

And Greek lowercase:

```

157   \um@mathcode@offset{#1}{945}{120572}%
158 }%
159 }

```

Uppercase Latin	<i>ABCDEFGHIJKLMNOPQRSTUVWXYZ</i>
Lowercase Latin	<i>abcdefghijklmnopqrstuvwxyz</i>
Uppercase Greek	<i>ΑΒΓΔΕΖΗΘΙΚΛΜΝΞΟΠΣΤΥΦΧΨ</i>
Lowercase Greek	<i>αβγδεζηθικλμνξοπρστυφχψ</i>

File II

STIX table data extraction

The source for the \TeX names for the very large number of mathematical glyphs are provided via Barbara Beeton's table file for the STIX project (ams.org/STIX). A version is located at `:// . . / / / - .` but it's not currently up to date.

A single file is produced containing all 3298 symbols. Future optimisations might include generating various (possibly overlapping) subsets so not all definitions must be read just to redefine a small range of symbols..

```

1 #!/bin/sh
2
3 cat stix-tbl.asc |
4 awk '
5 BEGIN {OFS="|"}
6 {if (usv != substr($0,2,5) )
7   {if (substr($0,2,1) != " ")
8     {usv = substr($0,2,5);
9     texname = substr($0,84,25);
10    type = substr($0,57,1);
11    description = tolower(substr($0,233,350));
12    {if (texname ~ /[\\]/)
13      print usv, texname, type, description;}}}' - |
14 awk -F"| " '
15 ((($3 != " ") && ($3 != "F") && ($3 != "D"))) {
16   print "\\unicode@math@symbol{" "\"" $1 "\"}-{\" $2 "\"}-{\" $3 "\"}-{\" $4 "\"};

```

```

17 }' - |
18 sed -e ' s/{N}/{\\mathord}/ ' \
19      -e ' s/{A}/{\\mathalpha}/ ' \
20      -e ' s/{P}/{\\mathpunct}/ ' \
21      -e ' s/{B}/{\\mathbin}/ ' \
22      -e ' s/{R}/{\\mathrel}/ ' \
23      -e ' s/{L}/{\\mathop}/ ' \
24      -e ' s/{O}/{\\mathopen}/ ' \
25      -e ' s/{C}/{\\mathclose}/ ' > stix-tex.tex

```