

UnicodeMath

A Nearly Plain-Text Encoding of Mathematics

Version 3.3

Murray Sargent III

22-Jan-25

1.	Introduction.....	2
2.	Encoding Simple Math Expressions.....	4
2.1	Fractions	4
2.2	Subscripts and Superscripts.....	6
2.3	Use of the Blank (Space) Character.....	8
3.	Encoding Other Math Expressions.....	9
3.1	Delimiters.....	9
3.2	Literal Operators	11
3.3	Prescripts and Above/Below Scripts.....	12
3.4	<i>n</i> -ary Operators	12
3.5	Mathematical Functions	13
3.6	Square Roots and Radicals	14
3.7	Enclosures.....	14
3.8	Stretchy Characters.....	16
3.9	Matrices.....	17
3.10	Accent Operators.....	18
3.11	Differential, Exponential, and Imaginary Symbols.....	19
3.12	Unicode Subscripts and Superscripts	19
3.13	Concatenation Operators.....	19
3.14	Comma, Period, and Colon.....	19
3.15	Ordinary Text Inside Math Zones	20
3.16	Space Characters	21
3.17	Phantoms and Smashes	22
3.18	Arbitrary Groupings.....	23
3.19	Equation Arrays.....	23
3.20	Math Zones.....	23
3.21	Equation Numbers	24
3.22	UnicodeMath Characters and Operands.....	24
3.23	Equation Breaking and Alignment.....	27
3.24	Size Overrides.....	27
3.25	Math Color	28
4.	Input Methods	28
4.1	Character Translations	29
4.2	Math Keyboards.....	30
4.3	Hexadecimal Input	31
4.4	Pull-Down Menus, Ribbons, Context Menus	31
4.5	Macros	31
4.6	UnicodeMath Autocorrect List	31
4.7	Handwritten Input	32

4.8	Math Dictation	32
4.9	Braille	33
5.	Recognizing Mathematical Expressions	33
5.1	Characters that may start a math zone.....	34
5.2	Characters that do start a math zone.....	34
5.3	Characters that confirm a possible math zone start.....	35
5.4	Spaces and vertical bars inside a math zone.....	35
5.5	Determining the math-zone end	35
5.6	Example.....	36
6.	Using UnicodeMath in Programming Languages.....	36
6.1	Advantages of UnicodeMath in Programs	37
6.2	Comparison of Programming Notations.....	38
6.3	Export to TeX.....	40
7.	Conclusions	41
	Acknowledgements	42
	Appendix A. UnicodeMath Grammar	42
	Appendix B. Character Keywords and Properties	43
	Version Differences	59
	References.....	60

1. Introduction

With a few conventions, [Unicode¹](#) can encode many mathematical expressions in readable nearly plain text. This format is a “lightly marked up format”; hence the use of “nearly”. The format is linear, but it can be displayed in built-up presentation form. We refer to the linear format as *UnicodeMath* and to the built-up presentation format as the *built-up format*. To see UnicodeMath in action, click on the Demo button in [UnicodeMathML](#), an open-source web applet that displays equations or expressions in UnicodeMath, in built-up format, and in MathML.

UnicodeMath can be used with heuristics based on the Unicode math properties to recognize mathematical expressions without the aid of explicit math-on/off commands. The recognition is facilitated by Unicode’s strong [support](#) for mathematical symbols.² Alternatively, UnicodeMath can be used in “math zones” explicitly controlled by the user either with explicit math zone delimiters as used in [La]TeX or with a character format attribute in a rich-text environment as in Microsoft Word. Use of math zones is desirable, since recognition heuristics are not infallible.

UnicodeMath is more compact and easy to read than [La]TeX,^{3,4} or MathML.⁵ However unlike those formats, it doesn’t attempt to include all typographical embellishments. Instead, some embellishments are delegated to the higher-level layer that handles rich-text properties like font size, footnotes, comments, hyperlinks, etc. In principle one can extend the notation to include the properties of the higher-level layer, but at the cost of reduced readability. Hence embedded in a rich-text environment, UnicodeMath can faithfully represent rich mathematical text, whereas embed-

ded in a plain-text environment it lacks most rich-text properties and some mathematical typographical properties. UnicodeMath is primarily concerned with presentation, but it has some semantic features that might seem to be only content oriented, e.g., n -aryands (integrands, summands, etc.) and function-apply arguments (see Secs. [3.4](#) and [3.5](#)). These aid in displaying built-up functions with proper typography and they also help to interoperate with math-oriented programs and math speech.

Most mathematical expressions can be represented unambiguously in UnicodeMath, from which they can be exported to [La]TeX, MathML, C++, speech, braille, and symbolic manipulation programs. UnicodeMath borrows notation from TeX for mathematical objects that don't lend themselves well to mathematical linear notation, e.g., for matrices. UnicodeMath input methods convert [La]TeX control words to the corresponding Unicode characters.

A variety of syntax choices can be used for a linear format. The choices made for UnicodeMath favor several criteria: efficient input of mathematical formulae, sufficient generality to support high-quality mathematical typography, the ability to round trip elegant mathematical text, at least in a rich-text environment, and a format that resembles mathematical notation. Some compromises between these goals had to be made.

UnicodeMath is useful for 1) inputting mathematical expressions,⁶ 2) displaying mathematics by text engines that cannot display a built-up format, and 3) computer programs. In addition to being the most readable linear math format, UnicodeMath is the most concise. It represents the simple fraction, one half, by the 3 characters “1/2”, whereas typical MathML takes 35 characters (`<mfrac><mn>1</mn><mn>2</mn></mfrac>`). This conciseness makes UnicodeMath an attractive format for storing mathematical expressions and equations, as well as for ease of keyboard entry. Another comparison is in the math structures for the Equation Tools tab in the Microsoft 365 math ribbon. In Word, the structures are defined in OMML (Office MathML) and built up by Word, while in the other math-enabled apps, the structures are defined in UnicodeMath and built up by RichEdit. The latter build up much faster and the equation data is much smaller. A dramatic example is the stacked-fraction template (empty numerator over empty denominator). In UnicodeMath, this is given by the single character ‘/’ while in OMML, it's 109 characters. LaTeX is considerably shorter at 9 characters “`\frac{}{}`”, but is still 9 times longer than UnicodeMath. [AsciiMath](#) represents fractions the same way as UnicodeMath, so simple cases are identical. If Greek letters or other characters that require names in AsciiMath are used, UnicodeMath is shorter and more readable.

Another advantage of UnicodeMath over MathML and OMML is that UnicodeMath can be stored anywhere Unicode text is stored. When adding math capabilities to a program, XML formats require redefining the program's file format and potentially destabilizing backward compatibility, while UnicodeMath does not. If a program is aware of UnicodeMath math zones (see Section [3.20](#)), it can recover the built-up mathematics by passing those zones through the RichEdit UnicodeMath [MathBuildUp](#)

function. In fact, you can roundtrip RichEdit documents containing math zones through plain-text editors like Notepad: the math zones are preserved.

For interchange of math expressions between arbitrary math-aware programs, MathML and other higher-level languages are preferred. Currently, conversion between UnicodeMath and other math formats is only implemented in Microsoft desktop apps, the Windows Calculator applet, and on the web in [UnicodeMathML](#), although there's nothing proprietary about UnicodeMath.

Section 2 motivates and illustrates UnicodeMath using the fraction, subscripts, and superscripts along with a discussion of how the ASCII space U+0020 is used to build up one construct at a time. [Section 3](#) summarizes the usage of the other constructs along with their relative precedences, which are used to simplify the notation. [Section 4](#) discusses input methods. [Section 5](#) gives an algorithm for recognizing mathematical expressions embedded in ordinary text. [Section 6](#) explains how Unicode plain text can be helpful in programming languages. [Section 7](#) gives conclusions. The appendices present a simplified [UnicodeMath grammar](#) and a list of [operators](#).

2. Encoding Simple Math Expressions

Given Unicode's strong support for mathematics² relative to ASCII, how much better can a plain-text encoding of mathematical expressions look using Unicode? The most well-known ASCII encoding of such expressions is that of TeX, so we use it for comparison. MathML is more verbose than TeX and some of the comparisons apply to it as well. Notwithstanding TeX's phenomenal success in the science and engineering communities, a casual glance at its representations of mathematical expressions reveals that they do not look very much like the expressions they represent. It's not easy to make algebraic calculations by hand using TeX's notation. With UnicodeMath, one can represent mathematical expressions more readably, and the results can often be used with few or no modifications for such calculations. This capability is considerably enhanced by using UnicodeMath in a system that can also display and edit the mathematics in built-up form, such as Microsoft 365 math-aware applications like Microsoft Word.

The present section introduces UnicodeMath with fractions, subscripts, and superscripts. It concludes with a subsection on how the ASCII space character U+0020 can be used to build up one construct at a time. This is a key idea that helps make UnicodeMath ideal for inputting mathematical formulae. In general, where syntax and semantic choices were made, input convenience was given higher priority.

2.1 Fractions

One way to specify a fraction linearly is LaTeX's `\frac{numerator}{denominator}`. The `{ }` are not printed when the fraction is built up. These simple rules immediately give a "plain text" that is unambiguous, but looks quite different from the corresponding mathematical notation, thereby making it harder to read.

Instead, we define a simple operand to consist of all consecutive letters and decimal digits, i.e., a span of alphanumeric characters, those belonging to the Lx and Nd General Categories (see [The Unicode Standard](#),¹ Table 4-2. General Category). As such, a simple numerator or denominator is terminated by most nonalphanumeric characters, including, for example, arithmetic operators, the blank (U+0020), and Unicode characters in the ranges U+2200..U+23FF, U+2500..U+27FF, and U+2900..U+2AFF. The fraction operator is given by the usual solidus / (U+002F). So the simple built-up fraction

$$\frac{abc}{d}$$

appears in UnicodeMath as `abc/d`. To force a display of a normal-size linear fraction, one can use `\/` (backslash followed by slash).

For more complicated operands (such as those that include binary or relational operators), parentheses (), brackets [], or braces { } can be used to enclose the desired character combinations. If parentheses are used and the outermost parentheses are preceded and followed by operators, those parentheses are not displayed in built-up form, since usually one does not want to see such parentheses. So, the plain text `(a + c)/d` displays as

$$\frac{a + c}{d}$$

In practice, this approach leads to plain text that is easier to read than LaTeX's, e.g., `\frac{a + c}{d}`, since in many cases, parentheses are not needed, while TeX requires `{}`'s. To force the display of the outermost parentheses, one encloses them, in turn, within parentheses, which then become the outermost parentheses. For example, `((a + c))/d` displays as

$$\frac{(a + c)}{d}$$

A neat feature of this notation is that the plain text is, in fact, often a legitimate mathematical notation in its own right, so it is relatively easy to read. Contrast this with the MathML version, which (with no parentheses) reads as

```
<mfrac>
  <mrow>
    <mi>a</mi>
    <mo>+</mo>
    <mi>c</mi>
  </mrow>
  <mi>d</mi>
</mfrac>
```

Three built-up fraction styles are available: the stacked fraction / U+002F, the “fraction slash” U+2044 (which one can input by typing `\sdiv`) which builds up to a

skewed fraction, and the “division slash” U+2215 (`\ldiv`) which builds up to a linear fraction. Examples are

$$\frac{\frac{a}{b+c}}{\frac{d}{e}+f}, \quad \frac{a}{b+c} / \frac{d}{e} + f, \quad \left(\frac{a}{b+c}\right) / \left(\frac{d}{e} + f\right)$$

When building up the large linear fraction, the outermost parentheses should not be removed.

The same notational syntax is used for a “stack” which is like a fraction with no fraction bar. The stack is used to create binomial coefficients and the stack operator is `'` (`\atop`). For example, the binomial theorem

$$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}$$

in UnicodeMath reads as

$$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k},$$

where $\binom{n}{k}$ is the binomial coefficient for the combinations of n items grouped k at a time. The summation limits use the subscript/superscript notation discussed in the next subsection. Since binomial coefficients are quite common, TeX has the `\choose` control word for them. In UnicodeMath, this uses the `\choose` operator (`^c`) instead of the `\atop` operator `'`. Accordingly, the binomial coefficient in the binomial theorem above can be written as “ $n \setminus \text{choose } k$ ”, if an operator follows the k .

When `/` is followed by an operator, it’s highly unlikely that a fraction is intended. This fact leads to a simple way to enter *negated* operators like \neq , namely, just type `/=` to get \neq . A list of such negated operator combinations is given in [Section 4.1](#). To enter \neq , you can also type TeX’s name, `\ne`, but `/=` is slightly simpler. And the TeX names for the other negated operators in Section 4.1 are harder to remember. One other trick with fractions is that a period or comma in between two digits or in between the slash and a digit is considered to be part of a number, rather than being a terminator. For example `1/3.1416` builds up to $\frac{1}{3.1416}$, rather than $\frac{1}{3}.1416$.

The fraction operators have left-to-right associativity as in common programming languages like C/C++/C#. For example, `1+a/b/c/d` builds up as

$$1 + \frac{\frac{a}{b}}{c/d}$$

2.2 Subscripts and Superscripts

Subscripts and superscripts are a bit trickier, but they’re still quite readable. Specifically, we introduce a subscript by a subscript operator, which we display as the ASCII underscore `_` as in TeX. A simple subscript operand consists of the string of one or more characters with the [Unicode General Categories](#) `Lx` (alphabetic) and `Nd` (decimal digits), as well as the invisible comma. For example, a pair of subscripts, such as

$\delta_{\mu\nu}$ is written as $\delta_{\mu\nu}$. Similarly, superscripts are introduced by a superscript operator, which we display as the ASCII `^` as in TeX. So, a^b means a^b . A nice enhancement for a text processing system with build-up capabilities is to display the `_` as a small subscript down arrow \downarrow and the `^` as a small superscript up arrow \uparrow , in order to convey the semantics of these build-up operators in a math context.

Compound subscripts and superscripts include expressions within parentheses, square brackets, and curly braces. So $\delta_{\mu+\nu}$ is written as $\delta_{(\mu+\nu)}$. In addition, it is worthwhile to treat two more operators, the comma and the period, in special ways. Specifically, if a subscript operand is followed directly by a comma or a period that is, in turn, followed by whitespace, then the comma or period appears on the baseline, i.e., is treated as the operator that terminates the subscript. However, a comma or period followed by an alphanumeric is treated as part of the subscript. This refinement obviates the need for many overriding parentheses, thereby yielding a more readable linear-format text (see [Sec. 3.14](#) for more discussion of comma and period).

Another kind of compound subscript is a subscripted subscript, which works using right-to-left associativity, e.g., a_{b_c} stands for a_{b_c} . Similarly, a^{b^c} stands for a^{b^c} . Fortran's `**` exponentiation operator also has right-to-left associativity.

Parentheses are needed for constructs such as a subscripted superscript like a^{b_c} , which is given by $a^{(b_c)}$, since a^b_c displays as a_c^b (as does a_c^b). The build-up program is responsible for figuring out what the subscript or superscript base is. Typically, the base is just a single math italic character like the a in these examples. But it could be a bracketed expression or the name of a mathematical function like \sin as in $\sin^2 x$, which renders as $\sin^2 x$ (see [Sec. 3.5](#) for more discussion of this case). It can also be an operator, as in the examples $+_1$ and $=_2$. In Indic and other cluster-oriented scripts the base is by default the cluster preceding the subscript or superscript operator.

As an example of a slightly more complicated example, consider the expression $W_{\delta_1\rho_1\sigma_2}^{3\beta}$, which can be written in UnicodeMath as $W^{3\beta}_{\delta_1\rho_1\sigma_2}$, where Unicode numeric subscripts are used. It can also be written as $W^{3\beta}_{\delta_1\rho_1\sigma_2}$. In TeX, it's

$$W^{3\beta}_{\delta_1\rho_1\sigma_2}$$

The TeX version looks simpler using Unicode for the symbols, namely $W^{3\beta}_{\delta_1\rho_1\sigma_2}$ or $W^{3\beta}_{\delta_1\rho_1\sigma_2}$, since Unicode has a full set of decimal subscripts and superscripts. As a practical matter, numeric subscripts are typically entered using an underscore and the number followed by a space or an operator, so the major simplification is that fewer brackets are needed.

For the ratio

$$\frac{\alpha_2^3}{\beta_2^3 + \gamma_2^3}$$

UnicodeMath can read as $\alpha_2^3/(\beta_2^3 + \gamma_2^3)$, while the standard TeX version reads as

$$\alpha_2^3 \over \beta_2^3 + \gamma_2^3$$

The UnicodeMath text is a legitimate mathematical expression, while the TeX version bears no resemblance to a mathematical expression.

TeX becomes cumbersome for longer equations such as

$$W_{\delta_1\rho_1\sigma_2}^{3\beta} = U_{\delta_1\rho_1}^{3\beta} + \frac{1}{8\pi^2} \int_{\alpha_1}^{\alpha_2} d\alpha'_2 \left[\frac{U_{\delta_1\rho_1}^{2\beta} - \alpha'_2 U_{\rho_1\sigma_2}^{1\beta}}{U_{\rho_1\sigma_2}^{0\beta}} \right]$$

A UnicodeMath version of this reads as

$$W_{\delta_1\rho_1\sigma_2}^{3\beta} = U_{\delta_1\rho_1}^{3\beta} + 1/8\pi^2 \int_{\alpha_1}^{\alpha_2} d\alpha'_2 [(U_{\delta_1\rho_1}^{2\beta} - \alpha'_2 U_{\rho_1\sigma_2}^{1\beta}) / U_{\rho_1\sigma_2}^{0\beta}]$$

while the standard TeX version reads as

$$\begin{aligned} & \mathbb{W}_{\{\delta_1\rho_1\sigma_2\}^{3\beta}} \\ &= U_{\{\delta_1\rho_1\}^{3\beta}} + \{1 \over 8\pi^2\} \\ & \int_{\{\alpha_1\}^{\{\alpha_2\}} d\alpha_2' \left[\right. \\ & \{U_{\{\delta_1\rho_1\}^{2\beta}} - \alpha_2' \\ & U_{\{\rho_1\sigma_2\}^{1\beta}} \over \\ & U_{\{\rho_1\sigma_2\}^{0\beta}} \} \left. \right] \end{aligned}$$

2.3 Use of the Blank (Space) Character

The ASCII space character U+0020 is rarely needed for explicit spacing of built-up text since the spacing around operators is provided automatically by the math display engine (Sec. 3.16 discusses this automatic spacing). However, the space character is very useful for delimiting the operands of UnicodeMath. When the space plays this role, it is eliminated upon build up. So, if you type `\alpha` followed by a space to get α , the space is eliminated when the α replaces the `\alpha`. Similarly, `a_1 b_2` builds up as a_1b_2 with no intervening space.

Another example is that a space following the denominator of a fraction is eliminated, since it causes the fraction to build up. If a space precedes the numerator of a fraction, the space is eliminated since it may be necessary to delimit the start of the numerator. Similarly if a space is used before a function-apply construct (see Sec. 3.5) or before above/below scripts (see Sec. 3.3), it is eliminated since it delimits the start of those constructs.

In a nested subscript/superscript expression, the space builds up one script at a time. For example, to build up $a^b c$ to a^{b^c} , two spaces are needed if spaces are used to build up. Binary and relational operators like `+` and `=` build up the whole expression, since the operands are terminated by such operators. For example, `a^b c+` builds up to $a^{b^c} +$.

In TeX, the space character is also used to delimit control words like `\alpha` and does not appear in built-up form. A difference between UnicodeMath's space-character usage and TeX's is that in TeX, spaces are invariably eliminated in built-up display,

whereas in UnicodeMath blanks that don't delimit operands or keywords result in spacing. Additional spacing characters are discussed in Sec. [3.16](#).

One displayed use for spaces is in overriding the algorithm that decides that an ambiguous unary/binary operator like + or – is unary. If followed by a space, the operator is considered to be binary and the space isn't displayed. Spaces are also used to obtain the correct spacing around comma, period, and colon in various contexts (see Sec. [3.14](#)).

3. Encoding Other Math Expressions

The previous section describes how UnicodeMath encodes fractions, subscripts and superscripts and gives a feel for that format. The current section describes how other mathematical constructs are encoded in UnicodeMath and ends with a more formal discussion of the syntax.

3.1 Delimiters

Brackets [], braces { }, and parentheses () represent themselves in UnicodeMath, and a word processing system capable of displaying built-up formulas should be able to enlarge them to fit around what's inside them. In general, we refer to such characters as *delimiters*. A delimited pair need not consist of the same kinds of delimiters. For example, it's fine to open with [and close with } and one sees this usage in some mathematical documents. The closing delimiter can have a subscript and/or a superscript. Delimiters are called *fences* in MathML.

These choices suffice for most cases of interest. But to allow for use of a delimiter without a matching delimiter and to overrule the open/close character of delimiters, the special keywords `\open` and `\close` can be used. These translate to the box-drawings characters `⏟` and `⏟`, respectively. We use box drawings characters for the open and close delimiters because they aren't likely to be used as mathematical characters and they are readily available in fonts. If used before any character that isn't a delimiter of the opposite sense, the open/close delimiter acts as an invisible delimiter, defining the corresponding end of a delimited expression. A common use of this is the “cases” equation, such as

$$|x| = \begin{cases} x & \text{if } x \geq 0 \\ -x & \text{if } x < 0 \end{cases}$$

which has the UnicodeMath “`|x| = {█ (&x" if "x ≥ 0@-&x" if "x < 0)⏟`” (see Sec. [3.19](#) for a discussion of the equation-array operator `█`).

Because the cases construct is common, TeX has the `\cases` control word for it. This is implemented in UnicodeMath with the `\cases` operator `Ⓞ`. With this the equation above can be written as “`|x| = Ⓞ(&x" if "x ≥ 0@-&x" if "x < 0)`”, which is still a little strange, but you don't have to type the opening curly brace or the `\close` character (`⏟`).

The open and close delimiters can be used to overrule the normal open and close character of delimiters as in the admittedly strange, but nevertheless sometimes used,

expression “]a + b[”, which has the UnicodeMath “ $\lrcorner a+b\lrcorner$ [”. Note that a blank following an open or close delimiter is “eaten”. This is to allow an open delimiter to be followed by a normal delimiter without interpreting the pair as a single delimiter. See also Sec. 3.18 on how to make arbitrary groupings. If a \lrcorner needs to be treated as an empty open delimiter when it appears before a delimiter like | or], follow the \lrcorner by a space to force the open-delimiter interpretation.

To suppress automatic sizing and to choose specific sizes, \lrcorner is followed by a digit ‘0’ – ‘4’ with the meanings in the following table

Digit	Meaning
0	Don’t grow
1	TeX <code>\big</code>
2	TeX <code>\Big</code>
3	TeX <code>\bigg</code>
4	TeX <code>\Bigg</code>

It’s rarely necessary to use explicit sizes if the display system can break equations within bracketed expressions.

The usage of open and close delimiters in UnicodeMath is admittedly a compromise between the explicit nature of TeX and the desire for a legitimate math notation, but the flexibility can be worth the compromise especially when interoperating with ordinarily built-up text such as in a WYSIWYG math system. TeX uses `\left` and `\right` for this purpose instead of `\open` and `\close`. We use the latter since they apply to right-to-left mathematics used in many Arabic locales as well as to the usual left-to-right mathematics.

Absolute values are represented by the ASCII vertical bar | (U+007C). The evenness of its count at any given bracket nesting level typically determines whether the vertical bar is a close |. Specifically, the first appearance is considered to be an open | (unless subscripted or superscripted), the next a close | (unless following an operator), the next an open |, and so forth.

Nested absolute values can be handled unambiguously by discarding the outermost parentheses within an absolute value. For example, the built-up expression $||x| - |y||$ can have the UnicodeMath $\lrcorner(|x|-|y|)\lrcorner$. Some cases, such as this one, can be parsed without the clarifying parentheses by noting that a vertical bar | directly following an operator is an open |. But the example $|a|b-c|d|$ needs the clarifying parentheses since it can be interpreted as either $(|a|b)-(c|d|)$ or $|a(|b-c)|d|$. The usual algorithm gives the former, so if one wants the latter without the inner parentheses, one can type $\lrcorner(a|b-c|d)\lrcorner$.

Another way to specify absolute values unambiguously is to use the `\abs` enclosure, which has the operator symbol ^(a). This was introduced for [math dictation](#) since one typically says, “the absolute value of ...”. That way you know the expression is an absolute value and not the cardinality of a set or a determinant. Note that a determinant identifies itself if the vertical bars contain a matrix. Alternatively, you can use the

`\det` or `\vmatrix` enclosures. Vertical bars are also used to denote cardinality of a set. To input this semantic, use `\card` which is represented in UnicodeMath by © (U+2402).

Another case where we treat `|` as a close delimiter is if it is followed by a space (U+0020). This handles the important case of the bra vector $\langle |$ in Dirac notation. For example, the quantum mechanical density operator ρ has the definition

$$\rho = \sum_{\psi} P_{\psi} |\psi\rangle\langle\psi|,$$

where the vertical bars can be input using the ASCII vertical bar.

If a `|` is followed by a subscript and/or a superscript and has no corresponding open `|`, it is treated as a script base character, i.e., *not* a delimiter. Its built-up size should be the height of the integral sign in the current display/inline mode.

The Unicode norm delimiter U+2016 (`\|` or `\norm`) has the same open/close definitions as the absolute value character `|` except that it's always considered to be a delimiter.

Delimiters can also have separators within them. UnicodeMath doesn't formalize the comma separators of function arguments (MathML does), but it supports the vertical bar separator `\vbar`, which is represented by the box drawings light vertical character `|` (U+2502). We tried using the ASCII `|` (U+007C) for this purpose too, but ambiguities can occur. A case using U+007C as a separator that's unambiguous has the form $(a|b)$, where a and b are mathematical expressions. But $(a|b|c)$ interprets the vertical bars as the absolute value. And one might want to interpret the `|` in $(a|b)$ as an open delimiter with `)` as the corresponding close delimiter, while the `(` isn't yet matched. If so, precede the `|` by `⋮`, i.e., $(\vdash b)$. The vertical bar separator grows in size to match the size of the surrounding brackets. Other operators can be treated as separators by preceding them with `\middle` (`\|`— U+2551).

Another common separator is the `\mid` character `|` (U+2223), commonly used in expressions like $\{x \mid f(x) = 0\}$. This separator also grows in size to match the surrounding brackets and is spaced as a relational operator.

3.2 Literal Operators

Certain operators like brackets, braces, parentheses, superscript, subscript, integral, etc., have special meaning in UnicodeMath. In fact, even a character like `+` plays a role in UnicodeMath in that it terminates an operand. To remove the UnicodeMath role of such an operator, we precede it by the “literal operator”, for which the backslash `\` is handy. So `\[` is displayed as an ordinary left square bracket, with no attempt by the build-up software to match a corresponding right square bracket. Such *quoted* operators are automatically included in the current operand.

All UnicodeMath operators consist of a single Unicode character, although a control word like `\open` may be used to input the character. Using a single character has the advantage of being globalized, since the control word typically looks like English. Users can define other control words that look like words in other languages just so

long as they map into the appropriate operator characters. A slight exception to the single-character operator rule occurs for accent operators (see Sec. 3.10). For these the accent combining mark may be preceded by a no-break space (U+00A0) for the sake of readability. Another advantage of using operator characters rather than control words is that the build-up processing is simplified and therefore faster. And it's helpful that operator characters look like the operators they represent, while the control words do not.

3.3 Prescripts and Above/Below Scripts

A special parenthesized syntax is used to form prescripts, that is, subscripts and superscripts that precede their base. For this $(_c^b)a$ creates the prescribed variable ${}^b{}_c a$. Variables can have both prescripts and postscripts (ordinary subscripts and superscripts).

You can use a prescript notation like TeX's. Just type a subscript and/or a superscript not preceded by a base and then follow it with a character that can be used as a base. For ${}^b{}_c a$, you type $_c^b a$. Note that you need to terminate the superscript with a space. If a variable precedes the prescript, you also need to precede the prescript with a space. A common use of prescripts is for the confluent hypergeometric functions, such as ${}_1F_1$. This can be input as $_1 F_1$ or as $(_1^)F_1$.

Below scripts and above scripts are represented in general by the line drawing operators $\backslash\text{below}$ (\mp) and $\backslash\text{above}$ (\pm), respectively. Hence the expression $\lim_{n \rightarrow \infty} a_n$ can be represented by $\lim_{\mp} (n \rightarrow \infty) a_n$. Since the operations \det , \gcd , \inf , \lim , \liminf , \limsup , \max , \min , \Pr , and \sup are common, their below scripts are also accessible by the usual subscript operator $_$. So, in display mode, $\lim_{n \rightarrow \infty} a_n$ can also be represented by $\lim_{_}(n \rightarrow \infty) a_n$, which is a little easier to type than $\lim_{\mp} (n \rightarrow \infty) a_n$.

Although for illustration purposes, the below script examples are shown here in-line with the script below, ordinarily this choice is only for display-mode math. When inline, below- and above-scripts entered with $_$ and $\^$ are shown as subscripts and superscripts, respectively, as are the limits for n -ary operators. When entered with \mp and \pm , they remain below and above scripts in-line. If an above/below operator or a subscript/superscript operator is preceded by an operator, that operator becomes the base. See Sec. 3.8 for some examples.

3.4 n -ary Operators

n -ary operators, sometimes called large operators, like integral, summation and product are sub/superscripted or above/below operators that have a third argument: the " n -aryand". For the integral, the n -aryand is the integrand, and for the summation, it's the summand. For both typographical and semantic purposes, it's useful to identify these n -aryands. In UnicodeMath, this is done by following the sub/superscripted n -ary operator by the naryand concatenation operator $\backslash\text{naryand}$ (⦿) which is U+2592, or by a space. The operand that follows this operator becomes the n -aryand. For example, the UnicodeMath expression $\int_0^a x dx / (x^2 + a^2)$ has the built up form

$$\int_0^a \frac{x dx}{x^2 + a^2}$$

where $x\mathit{d}x/(x^2+a^2)$ is the integrand and d is the Unicode differential d character U+2146. Unlike for the fraction numerator and denominator, the outermost parentheses of an n -aryand are *not* removed on buildup, since parentheses are commonly used to delimit compound n -aryands. Notice that the d character automatically introduces a small space between the x and the dx and by default displays as a math-italic d when it appears in a math zone (this choice could be changed by a document default setting).

To delimit more complicated n -aryands without using parentheses or brackets of some kind, use the `\begin \end ([]` see Sec. 3.18) delimiters, which disappear on build up.

Since `\naryand` isn't the most intuitive name, the alias `\of` can be used. This also works as an alias for `\funcapply` in math function contexts (see Sec. 3.5). This alias is motivated by sentences like "The integral from 0 to b of $x dx$ is one-half b squared."

Sometimes one wants to control the positions of the limit expressions explicitly as in using TeX's `\limits` (upper limit above, lower below) and `\nolimits` (upper limit as superscript and lower as subscript) control words. To this end, if the n -ary operator is followed by the digit 1, the limit expressions are displayed above and below the n -ary operator and if followed by the digit 2, they are displayed as superscript and subscript. More completely, the number can be one of the first five of the following, OR'd with any of the next three, along with neither or one of the last two.

<code>nLimitsDefault</code>	0
<code>nLimitsUnderOver</code>	1
<code>nLimitsSubSup</code>	2
<code>nUpperLimitAsSuperScript</code>	3
<code>nLimitsOpposite</code>	4
<code>nShowLowLimitPlaceHolder</code>	8
<code>nShowUpLimitPlaceHolder</code>	16
<code>fDontGrowWithContent</code>	64
<code>fGrowWithContent</code>	128

3.5 Mathematical Functions

Mathematical functions such as trigonometric functions like "sin" should be recognized as such and not italicized. As such they are treated as ordinary text (see Sec. 3.16). In addition, it's desirable to follow them with the Invisible Function Apply operator U+2061 (`\funcapply`). This is a special binary operator and the operand that

follows it is the function argument. In converting to built-up form, this operator transforms its operands into a two-argument object that renders with the proper spacing for mathematical functions.

If the Function Apply operator is immediately followed by a subscript or superscript expression, that expression should be applied to the function name and the Function Apply operator moved passed the modified name to bind the operand that follows as the function argument. For example, the function $\sin^2 x$ falls into this category.

Unlike for the fraction numerator and denominator, the outermost parentheses of the second operand of the function-apply operator are not removed on buildup, since parentheses are commonly used to delimit function arguments. To delimit a more complicated arguments without using parentheses or brackets of some kind, use the $\llbracket \rrbracket$ delimiters (`\begin \end`) which disappear on build up. If brackets are used, they and their included content comprise the function’s argument. For example, $\sin(x) b$ means $\sin(x) \times b$. To get $\sin(\omega - \omega_0)t$, where t is part of the argument, enclose the argument in $\llbracket \rrbracket$ delimiters.

Since `\funcapply` isn’t the most intuitive name, `\of` can be used in function-apply contexts. `\of` autocorrects to 𐀀 (U+2592—`\naryand`, see Sec. 3.4), but context can give it this convenient second use. This alias is motivated by sentences like “The sine of $2x$ equals twice the sine of x times the cosine of x ”, i.e., $\sin 2x = 2 \sin x \cos x$.

If a function name has a space in it, e.g., “lim sup”, the space is represented by a no-break space (U+00A0) as described in Sec. 3.16. If an ordinary ASCII space is used, it builds up the “lim” function.

3.6 Square Roots and Radicals

Square, cube, and quartic roots can be represented by expressions started by the corresponding Unicode radical characters $\sqrt{}$ (U+221A, `\sqrt`), $\sqrt[3]{}$ (U+221B, `\cbrt`), and $\sqrt[4]{}$ (U+221C, `\qdr`). These operators include the operand that follows. Examples are \sqrt{abc} , $\sqrt{a+b}$ and $\sqrt[3]{c+d}$, which display as \sqrt{abc} , $\sqrt{a+b}$, and $\sqrt[3]{c+d}$, respectively. In general, the n th root radical is represented by an expression like $\sqrt[n]{a}$, where a is the complete radicand. Anything following the closing parenthesis is not part of the radicand. For example, $\sqrt[n]{a+b}$ displays as $\sqrt[n]{a+b}$.

You can also obtain $\sqrt[n]{a+b}$ using the more TeX-like input `\root n\of(a+b)`. In this format, the degree of the radical can be more than one character without enclosing it in parentheses. For example, $\sqrt[n+1]{b+c}$ can be input by `\root n+1\of(b+c)`, which is like TeX’s `\root n+1\of{b+c}`.

3.7 Enclosures

To enclose an expression in a rectangle, use the rectangle operator $\boxed{}$ (U+25AD, `\rect`) followed by the operand representing the expression. This syntax is similar to that for the square root. For example, $\boxed{E = mc^2}$ displays as $\boxed{E = mc^2}$. Use the

same approach to put an overbar above an expression, namely follow the overbar operator $\overline{}$ (U+00AF, `\overbar`) by the desired operand. For an underbar, use the operator $\underline{}$ (U+2581, `\underbar`).

In general, the rectangle function can represent any combination of borders, horizontal, vertical, and diagonal strikeouts, and enclosure forms defined by the MathML `<menclase>` element, except for roots, which are represented as discussed in the previous Section. The general syntax for enclosing an expression x is $\boxed{n&x}$, where n is a mask consisting of any combination of the following flags:

fBoxHideTop	1
fBoxHideBottom	2
fBoxHideLeft	4
fBoxHideRight	8
fBoxStrikeH	16
fBoxStrikeV	32
fBoxStrikeTLBR	64
fBoxStrikeBLTR	128

It is anticipated that the enclosure format number n is chosen via some kind of friendly user interface, but at least the choice can be preserved in UnicodeMath. Note that the overbar function (`\overbar` U+00AF) can also be given by $\boxed{14&x}$, which hides all but the top border. The underbar (`\underbar` U+2581) can also be given by $\boxed{13&x}$, which hides all but the bottom border.

Other enclosures such as rounded box $\boxed{}$, circle \bigcirc , long division $\overline{}$, actuarial, and ellipse \ominus can be encoded as for the rectangle operator but using appropriate Unicode characters (not all chosen here).

An abstract box can be put around an expression x to change alignment, spacing category, size style, and other properties. This is defined by $\boxed{n&x}$, where $\boxed{}$ is U+25A1 (`\box`) and n can be a combination of one Align option, one Space option, one Size option and any flags in the following table:

nAlignBaseline	0
nAlignCenter	1
nSpaceDefault	0
nSpaceUnary	4
nSpaceBinary	8
nSpaceRelational	12
nSpaceSkip	16
nSpaceOrd	20
nSpaceDifferential	24
nSizeDefault	0
nSizeText	32

nSizeScript	64
nSizeScriptScript	96
fBreakable	128
fXPositioning	256
fXSpacing	512

3.8 Stretchy Characters

In addition to overbars and underbars, stretchable brackets are used in mathematical text. For example, the “underbrace” and “overbrace” are as

$$\begin{array}{c} k \text{ times} \\ \overbrace{x + \cdots + x} \\ \underbrace{x + y + z} \\ >0 \end{array}$$

The UnicodeMath for these are $\overbrace{(x+\cdots+x)}^{k \text{ "times"}}$ and $\underbrace{(x+y+z)}_{>0}$, respectively. Here the subscript and superscript operators are used for convenient keyboard entry (and compatibility with TeX); one can also use Sec. 3.3’s below-script and above-script operators, respectively. The horizontal stretchable brackets are given in the following table

U+23DC	$\overparen{}$	<code>\overparen</code>
U+23DD	$\underparen{}$	<code>\underparen</code>
U+23DE	$\overbracket{}$	<code>\overbrace</code>
U+23DF	$\underbracket{}$	<code>\underbrace</code>
U+23E0	$\overshell{}$	<code>\overshell</code>
U+23E1	$\undershell{}$	<code>\undershell</code>
U+23B4	$\overbracket{}$	<code>\overbracket</code>
U+23B5	$\underbracket{}$	<code>\underbracket</code>

There are many other characters that can stretch horizontally to fit text, such as various horizontal arrows. There are four configurations: a stretch character above or below a baseline text, and text above or below a baseline stretched character. Illustrating UnicodeMath for these four cases with the stretchy character \rightarrow and the text $a + b$, we have

$(a + b)^{\rightarrow}$	$\overrightarrow{a + b}$
$(a + b)_{\rightarrow}$	$\underline{\rightarrow} a + b$
$\rightarrow^{\perp} a + b$	$\overrightarrow{a + b}$
$\rightarrow_{\perp} (a + b)$	$\underline{\rightarrow} a + b$

3.9 Matrices

Matrices are represented by a notation similar to TeX's, namely an expression of the form

$$\blacksquare (exp_1 [\& exp_2] \dots @ \dots exp_{n-1} [\& exp_n] \dots)$$

where \blacksquare is the matrix character U+25A0 and @ is used to terminate rows, except for the last row which is terminated by the closing paren. This causes exp_1 to be aligned over exp_{n-1} , etc., to build up an $n \times m$ matrix array, where n is the maximum number of elements in a row and m is the number of rows. The matrix is constructed with enough columns to accommodate the row with the largest number of entries, with rows having fewer entries given sufficient null entries to keep the table $n \times m$. As an example, $\blacksquare(a\&b@c\&d)$ displays as

$$\begin{array}{cc} a & b \\ c & d \end{array}$$

If you want parentheses around the matrix, include them as in $(\blacksquare(a\&b@c\&d))$. Because parenthesized matrices are quite common, TeX has the `\pmatrix` control word that automatically includes parentheses. This is implemented with the `\pmatrix` operator $(^m)$. So $(^m(a\&b@c\&d))$ displays as

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

There are five matrix constructs with enclosing brackets that can be entered as summarized in the following table in which ... stands for the matrix contents.

LaTeX	Char	Code	Form
<code>\matrix</code>	\blacksquare	U+25A0	...
<code>\bmatrix</code>	$\textcircled{\text{S}}$	U+24E2	[...]
<code>\pmatrix</code>	$(^m)$	U+24A8	(...)
<code>\vmatrix</code>	$(^v)$	U+24B1	\dots
<code>\Bmatrix</code>	$\textcircled{\text{S}}$	U+24C8	\{\dots\}
<code>\Vmatrix</code>	$(^n)$	U+24A9	\ \dots\

Sometimes you just want to enter a sample matrix quickly. If any of the six matrix control words are followed by a digit d , they insert a $d \times d$ identity matrix. For example, typing `\pmatrix 3` inserts

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

This is easier to type than `\pmatrix(1&0&0@0&1&0@0&0&1)`, which displays the same identity matrix. An empty $n \times m$ matrix can be inserted by $n \times m$ <matrix control word>. For example, 2×3 `\pmatrix` not followed by ‘(‘ inserts

$$\left(\quad \right)$$

3.10 Accent Operators

Mathematics often has accented characters. Simple primed characters like a' are represented by the character followed by the Unicode prime U+2032, which can be typed in using the ASCII apostrophe '. Double primed characters have two Unicode primes, etc. In addition, Unicode has multiple prime characters that render with somewhat different spacing than concatenations of U+2032. The primes are special in that they need to be superscripted with appropriate use of heavier glyph variants (see Sec. 3.12). When a prime follows a variable, e.g., a' , it should be converted into a superscript function with a as the base and the prime as the superscript. It's also important to merge the prime into a superscript that follows, e.g., a'^c should display as a'^c , where both the prime and the c are in the same superscript argument.

The ASCII asterisk is raised in ordinary text, but in a math zone it gets translated into U+2217, which is placed on the math axis as the +. To make it a superscript or subscript, the user has to include it in a superscript or subscript expression. For example, a^{*2} has the UnicodeMath `a^*2` or `a^{*2}`. Here for convenience, the asterisk is treated as an operand character if it follows a subscript or superscript operator.

Other kinds of accented characters can be represented by Unicode combining mark sequences. The combining marks are found in the Unicode ranges U+0300—U+036F and U+20D0 – U+20FF. The most common math accents are summarized in the following table

<code>\hat</code>	U+0302	\hat{a}
<code>\check</code>	U+030C	\check{a}
<code>\tilde</code>	U+0303	\tilde{a}
<code>\acute</code>	U+0301	\acute{a}
<code>\grave</code>	U+0300	\grave{a}
<code>\dot</code>	U+0307	\dot{a}
<code>\ddot</code>	U+0308	\ddot{a}
<code>\dddotted</code>	U+20DB	\dddot{a}
<code>\bar</code>	U+0304	\bar{a}
<code>\vec</code>	U+20D7	\vec{a}

If a combining mark should be applied to more than one character or to an expression, that character or expression should be enclosed in parentheses and followed by the combining mark. Since this construct looks funny when rendered by

plain-text programs, a no-break space (U+00A0) can appear in between the parentheses and the combining mark. For example, $(a + b)^{\wedge}$ renders as $\widehat{a + b}$ when built up. Special cases of this notation include overscoring (use U+0305) and underscoring (use U+0332) mathematical expressions.

The combining marks are treated by a mathematics renderer as operators that translate into special accent built-up functions with the proper spacing for mathematical variables.

3.11 Differential, Exponential, and Imaginary Symbols

Unicode contains a number of special double-struck math italic symbols that are useful for both typographical and semantic purposes. These are U+2145—U+2149 for double-struck $D, d, e, i,$ and j ($\mathcal{D}, \mathcal{d}, \mathcal{e}, \mathcal{i}, \mathcal{j}$), respectively. They have the meanings of differential, differential, natural exponent, imaginary unit, and imaginary unit, respectively. They can be typed in using `\Dd`, `\dd`, `\ee`, `\ii`, and `\jj`, respectively.

In US patent applications these characters should be rendered as $\mathcal{D}, \mathcal{d}, \mathcal{e}, \mathcal{i}, \mathcal{j}$ as defined, but in regular US technical publications, these quantities can be rendered as math italic. In European technical publications, they are sometimes rendered as upright characters. Furthermore, the D and d start a differential expression and should have appropriate spacing for differentials. UnicodeMath treats these symbols as operand characters, but the display routines should provide the appropriate glyphs and spacings. See Sec. 3.4 for an example of an integral using \mathcal{d} .

3.12 Unicode Subscripts and Superscripts

Unicode contains a small set of mostly numeric superscripts (U+00B2, U+00B3, U+00B9, U+2070—U+207F) and a similar set of subscripts (U+2080—U+208F) that should be rendered the same way that scripts of the corresponding script nesting level would be rendered. To perform this translation, these characters can be treated as high-precedence operators, spans of which combine into the corresponding superscripts or subscripts when built up. Since numeric subscripts and superscripts are very common in mathematics, it's worthwhile building up Unicode subscripts and superscripts as if they had been UnicodeMath subscripts and superscripts.

3.13 Concatenation Operators

Relational and binary operators are “concatenation operators” so named because they are concatenated with their surrounding text in built-up form. In addition, a concatenation operator has two effects: 1) it terminates whatever operand precedes it, and 2) it implies appropriate surrounding space as discussed in Sec. 3.16 along with the mathematical spacing tables of the font. Since the spacing around operators is well-defined in this way, the user rarely needs to add explicit space characters.

3.14 Comma, Period, and Colon

The comma, period, and colon have context sensitive spacing requirements that can be represented in UnicodeMath.

Comma: when surrounded by ASCII digits render with ordinary text spacing. Else treat as punctuation with or without an ASCII blank following it. In either punctuation case the comma is displayed with a small space following it. If two spaces follow, the comma is rendered as a clause separator (a relatively large space follows the comma).

Period: when surrounded by ASCII digits render with ordinary text spacing. Else treat as punctuation with or without an ASCII blank following it. In either punctuation case the period is displayed with a small space following it. No clause separator option exists for the period. An extended decimal-point heuristic useful in calculator scenarios allows one to omit a leading 0, e.g., use numbers like .5. For this if the period is followed by an ASCII digit and 1) is at the start of a math zone, 2) follows a built-up math object start character or end-of-argument character, or 3) follows any operator except for closers and punctuation, then the period should be classified as a decimal point. With this algorithm, $a/.3$ displays as

$$\frac{a}{.3}$$

Colon: `<space> ‘:’` is displayed as Unicode RATIO U+2236 with relational spacing. `‘:’` without a leading space is displayed as itself with punctuation spacing.

3.15 Ordinary Text Inside Math Zones

Sometimes one wants ordinary text inside a function argument or in a math zone as in the formula

$$\text{rate} = \frac{\text{distance}}{\text{time}}$$

For such cases, the alphabetic characters should not be converted to math alphabetic characters and the typography should be that of ordinary text, not math text. To embed such text inside functions or in general in a math zone, the text can be enclosed inside ASCII double quotes. So, in UnicodeMath the formula above reads as

`"rate"="distance"/"time"`

If you want to include a double quote inside such text, insert `\"`. Another example is $\sin \theta = \frac{1}{2}e^{i\theta} + \text{c.c.}$ To get the “c.c.” as ordinary text, enclose it with ASCII double quotes. Otherwise the *c*’s will be italicized and the periods will have some space after them.

Alternatively ordinary text inside a math zone can be specified using a character-format property. This property is exported to plain text started and ended with the ASCII double quote. Note that no math object or math text can be nested inside an ordinary text region. Instead, if you paste a math object or text into an ordinary text

region, you split the region into two such regions with the math object and/or text in between.

3.16 Space Characters

Unicode contains numerous space characters with various widths and properties. These characters can be useful in tweaking the spacing in mathematical expressions. Unlike the ASCII space, which is removed when causing build up as discussed in Sec. 2.3, the other spaces are not removed on build up. Spaces of interest include the no-break space (U+00A0) and the spaces U+2000—U+200B, 202F, 205F.

In mathematical typography, the widths of spaces are usually given in integer multiples of an eighteenth of an em. The em space is given by U+2003. Various space widths are defined in the following table, which includes the corresponding MathML names having these widths by default

Space	Unicode	MathML name	Autocorrect
0 em	U+200B	zero-width space	<code>\zwsp</code>
1/18 em	U+200A	veryverythinmathspace	<code>\hairsp</code>
2/18 em	U+200A U+200A	verythinmathspace	
3/18 em	U+2009	thinmathspace	<code>\thinsp</code>
4/18 em	U+205F	mediummathspace	<code>\medsp</code>
5/18 em	U+2005	thickmathspace	<code>\thicksp</code>
6/18 em	U+2004	verythickmathspace	<code>\vthicksp</code>
7/18 em	U+2004 U+200A	veryverythickmathspace	
9/18 em	U+2002	ensp	<code>\ensp</code>
18/18 em	U+2003	emsp	<code>\emsp</code>
digit width	U+2007	numsp	<code>\numsp</code>
space width	U+00A0	no-break space	<code>\nbspsp</code>

In general, spaces act as concatenation operators and build up higher-precedence operators that precede them. But it's useful for the zero-width space (U+200B) to be treated as an operand character and not build up the preceding operator. The no-break space (U+00A0) is used when two words need to be separated by a blank but remain on the same line together. The no-break space is also treated as an operand character so that UnicodeMath combinations like “lim sup” and “lim inf” can be recognized as single operands. If an ASCII space (U+0020) appears after the “lim”, it builds up the “lim” function, rather than being part of the “lim sup” or “lim inf” function.

In math zones, most spacing is implied by the properties of the characters. The following table shows examples of how many 1/18^{ths} of an em size are automatically inserted between a character with the row property followed by a character with the column property for text-level expressions (see also p. 170 of *The TeXbook*)

	ord	unary	binary	rel	open	close	punct
ord	0	0	4	5	0	0	0
unary	0	0	4	0	0	0	0
binary	4	4	0	0	4	0	0
rel	5	5	0	0	5	0	0
open	0	0	0	0	0	0	0
close	0	0	4	5	0	0	0
punct	3	3	0	3	3	3	3

For the combinations described by this simple table, all script-level spacings are 0, but a more complete table would have some nonzero values. For example, in the expression $a + b$, the letters a and b have the ord (ordinary) property, while the $+$ has the binary property in this context. Accordingly, for the text level there is 4/18th em between the a and the $+$ and between the $+$ and the b . Similarly, there is 5/18th em between the $=$ and the surrounding letters in the equation $a = b$. A more complete table could include properties like math functions (trigonometric functions, etc.), n -ary operators, tall delimiters, differentials, subformulas (e.g., expression with an over brace), binary with no spacing (e.g., $/$), clause separators, ellipsis, factorial, and invisible function apply.

The zero-width space (U+200B, `\zwsp`) can be inserted into a null argument to suppress display of the argument place holder `⋮`. It can also be used to suppress subscript kerning. For example, the expression \mathcal{V}_{ab} shows the subscript ab automatically kerned in under the overhang of the \mathcal{V} . To prevent this kerning, one can insert a `\zwsp` before the subscript, which then displays unknerned as \mathcal{V}_{ab} .

3.17 Phantoms and Smashes

Sometimes one wants to obtain horizontal and/or vertical spacings that differ from the normal values. In [La]TeX this can be accomplished using phantoms to introduce extra space or smashes to zero out space. In UnicodeMath, seven special cases are defined as in the following table

Autocorrect	LF op	Op name	width	ascent	descent	ink
<code>\phantom</code>	◇ U+27E1	white concave-sided diamond	w	a	d	no
<code>\hphantom</code>	↔ U+2B04	white left-right arrow	w	0	0	no
<code>\vphantom</code>	↕ U+21F3	white up-down arrow	0	a	d	no
<code>\smash</code>	↓ U+2B0D	black up-down arrow	w	0	0	yes
<code>\asmash</code>	↑ U+2B06	black up arrow	w	0	d	yes
<code>\dsmash</code>	↓ U+2B07	black down arrow	w	a	0	yes
<code>\hsmash</code>	↔ U+2B0C	black left-right arrow	0	a	d	yes

The general case is given by `\phantom(n &<operand>)`, where n is any combination of the following flags:

fPhantomShow	1
fPhantomZeroWidth	2
fPhantomZeroAscent	4
fPhantomZeroDescent	8
fPhantomTransparent	16

For example, in the following equation the π in the upper limit is inside an `\hsmashphantom`, so that it has no width and thereby pulls the integrand in toward the integral sign

$$\frac{1}{2\pi} \int_0^{2\pi} \frac{d\theta}{a + b \sin \theta} = \frac{1}{\sqrt{a^2 - b^2}}$$

3.18 Arbitrary Groupings

The left/right white lenticular brackets `[` and `]` (U+3016 and U+3017) can be used to delimit an arbitrary expression without displaying these brackets on build up. The elimination of outermost parentheses for arguments of fractions, subscripts, and superscripts solves such grouping problems nicely in most cases, but the white lenticular brackets can handle any remaining cases. Note that in math zones, these brackets should be displayed using a math font rather than an East Asian font.

3.19 Equation Arrays

To align one equation relative to another vertically, one can use an equation array, such as

$$\begin{array}{l} 10x + 3y = 2 \\ 3x + 13y = 4 \end{array}$$

which has the UnicodeMath `[10&x+&3&y=2@3&x+&13&y=4]`, where `[` is U+2588 (input `\eqarray`). Here the meaning of the ampersands alternates between *align* and *spacer*, with an implied *spacer* at the start of the line. So, every odd `&` in a line is an alignment point and every even `&` is a place where space may be added to align the equations. This convention is used in [AMS-LaTeX](#).

3.20 Math Zones

[Section 5](#) gives an algorithm to identify the start and end of math zones in plain text. While the approaches given are surprisingly successful, they are not infallible. Hence if one knows the start and end of math zones, it's desirable to preserve this information in UnicodeMath.

In plain text, UnicodeMath uses `[` (U+2045) to start a math zone and `]` (U+2046) to end it. These characters are not ordinarily used in technical documents, so they would rarely need to be quoted (preceded by a backslash). When importing plain text, the user can execute a command to build up math zones defined by these math-zone delimiters.

The delimiters are analogous to TeX's `...` (inline math zones) and `$$...$$` (display math zones). UnicodeMath has the convention that if a math zone fills a (hard or soft) paragraph, the math zone is a display math zone. If any part of the paragraph

isn't in a math zone including a possible terminating period or comma, then the math zone is an inline math zone, which has more compact rendering. Adjacent math zones are automatically merged into a single math zone. Accordingly, UnicodeMath only needs one set of math zone delimiters. LaTeX display math zones can have the form $\left[\dots \right]$ and LaTeX inline math zones can have the form $\left(\dots \right)$.

3.21 Equation Numbers

Equation numbers are often used with equations presented in display mode. To represent an equation number flushed right of the equation in UnicodeMath, enter the equation followed by a # (U+0023) followed by the desired equation number text. For example $\left(E=mc^2 \right) \#(30)$ or more simply just $E=mc^2 \#(30)$ renders as

$$E = mc^2 \tag{30}$$

3.22 UnicodeMath Characters and Operands

UnicodeMath divides the ~150,000 assigned Unicode characters into three broad categories: 1) operand characters such as alphanumerics, 2) the bracket characters described in Sec. 3.1, and 3) other operator characters such as those described in Secs. 2.1–2.2 and 3.2–3.19. Operand characters include some nonalphanumeric characters, such as infinity (∞), exclamation point (!) if preceded by an operand, math minus (U+2212) or plus if either starts a sub/superscript operand, and period and comma if they're surrounded by ASCII (or full-width ASCII) digits (Sec. 3.14 gives a generalization of this last case). In other contexts, period and comma are treated as operators with the same precedence as plus. To reveal which characters are operators, operator-aware editors could be instructed to display operators with a different color or some other attribute.

In addition, operands include bracketed expressions and mixtures of such expressions and other operand characters. Hence $f(x)$ can be an operand. More specific definitions of operands are given in the simplified UnicodeMath syntax of [Appendix A](#).

Operands in subscripts, superscripts, fractions, roots, boxes, etc. are defined in part in terms of operators and operator precedence. While such notions are very familiar to mathematically oriented people, some of the symbols that we define as operators might be surprising. Most notably, the space (U+0020) is an important operator in UnicodeMath since it can be used to terminate operands as discussed in Sec. 2.3. A small but common list of operators is given in Table 3.1

Table 3.1 List of the most common operators ordered by increasing precedence

CR
([{ † [
)] } †]

&
space “ . , = - + * × · • ● ▫
/
∫ ∑ ∏
^
□ □ ■ √ ∛ ∜ − −
combining marks

where CR = U+000D. Note that the ASCII vertical bar | (U+007C) shows up both as an opening bracket and as a closing bracket. The choice is disambiguated by the evenness of its count at any given bracket nesting level or other considerations (see Sec. 3.1). So typically the first appearance is considered to be an open |, the next a close |, the next an open |, and so forth. The vertical bar appearing on the same level as & is considered to be a vertical bar separator and is given by the box drawings light vertical character (U+2502). We tried using the ASCII U+007C for this too, but the resulting ambiguities were insurmountable except in simple cases like (a|b) (see Section 3.1).

As in arithmetic, operators have precedence, which streamlines the interpretation of operands. The operators are grouped above in order of increasing precedence, with equal precedence values on the same line. For example, in arithmetic, $3+1/2 = 3.5$, not 2. Similarly, the UnicodeMath expression $\alpha + \beta/\gamma$ means

$$\alpha + \frac{\beta}{\gamma} \text{ not } \frac{\alpha + \beta}{\gamma}$$

Precedence can be overruled using parentheses, so $(\alpha + \beta)/\gamma$ gives the latter.

The following gives a list of the syntax for a variety of mathematical constructs (see Appendix A for a more complete grammar).

exp_1/exp_2	Create a built-up fraction with numerator exp_1 and denominator exp_2 . Numerator and denominator expressions are terminated by operators such as /*]) and blank (can be overruled by enclosing in parentheses).
$exp_1 exp_2$	Similar to fraction, but no fraction bar is displayed. Sometimes called a stack.
$base^exp_1$	Superscript expression exp_1 to the base $base$. The superscripts $^{0-9+-()}$ exist as Unicode symbols. Sub/superscript expressions are terminated, for example, by /*]) and blank. Sub/superscript operators associate right to left.
$base_exp_1$	Subscript expression exp_1 to the base $base$. The subscripts $_{0-9+-()}$ exist as Unicode symbols.

$base_exp_1^{exp_2}$	Subscript expression exp_1 and superscript expression exp_2 to the base $base$. The subscripts $_0$ – $_9$ + $_{-}$ $()$ exist as Unicode symbols.
$(_exp_1^{exp_2})base$	Prescript the subscript exp_1 and superscript exp_2 to the base $base$.
$base^{\perp}exp_1$	Display expression exp_1 centered above the base $base$. Above/below script operators associate right to left.
$base_{\top}exp_1$	Display expression exp_1 centered below the base $base$.
$[exp_1]$	Surround exp_1 with built-up brackets. Similarly for $\{ \}$ and $()$. Similarly for $\{ \}$, $()$, $ $. See Sec. 3.1 for generalizations.
$[exp_1]^{exp_2}$	Surround exp_1 with built-up brackets followed by superscripted exp_2 (moved up high enough).
$\square exp_1$	Abstract box around exp_1 .
$\square exp_1$	Rectangle around exp_1 .
$_exp_1$	Underbar under exp_1 (underbar operator is U+2581, not the ASCII underline character U+005F).
\bar{exp}_1	Overbar above exp_1 .
$\sqrt{exp_1}$	Square root of exp_1 .
$\sqrt[3]{exp_1}$	Cube root of exp_1 .
$\sqrt[4]{exp_1}$	Fourth root of exp_1 .
$\sqrt[exp_1]{exp_2}$	exp_1^{th} root of exp_2 .
$\Sigma_exp_1^{exp_2} exp_3$	Summation from exp_1 to exp_2 with summand exp_3 . $_exp_1$ and exp_2 are optional.
$\prod_exp_1^{exp_2} exp_3$	Product from exp_1 to exp_2 with multiplicand exp_3 . $_exp_1$ and exp_2 are optional.
$\int_exp_1^{exp_2} exp_3$	Integral from exp_1 to exp_2 with integrand exp_3 . $_exp_1$ and exp_2 are optional.
$(exp_1 [\& exp_2] \dots [exp_{n-1} [\& exp_n] \dots])$	Align exp_1 over exp_{n-1} , etc., to build up an array (see Appendix A for complete syntax).

Note that Unicode's plethora of mathematical operators² fill out the capabilities of the approach in representing mathematical expressions in UnicodeMath.

Precedence simplifies the text representing formulas but may need to be overruled. To terminate an operand (shown above as, for example, exp_1) that would otherwise combine with the following operand, insert a space (U+0020). This space does not show up when the expression is built up. Spaces that don't terminate operands may be used to space formulas in addition to the built-in spacing provided by a math display engine. Spaces are discussed in greater detail in Sec. [2-3](#).

To form a compound operand, parentheses can be used as described for the fraction above. For such operands, the outermost parentheses are removed. These operands occur for fraction numerators and denominators, subscript and superscript expressions, and arguments of functions like square root. Parentheses appearing in other contexts are always displayed in built-up format.

A curious aspect of the notation is that implied multiplication by juxtaposing two variable letters has very high precedence (just below that of diacritics), while explicit multiplication by asterisk and raised dot has a precedence equal to that of plus. So even though the analysis is like that for arithmetic expressions, it differs occasionally from the latter.

3.23 Equation Breaking and Alignment

Two features aid equation breaking and alignment in display math zones. A soft (optional) line break is created by the invisible times `\itimes` (U+2062), which is a binary operator, and you can break on it and align to it. It shouldn't display a glyph, except for a thin space if at the end of a math zone or if it's at the start of a line. With it you can effectively break an equation before any character, not just on binary, relational and some other operators. Generally, it's nice to display a multiplication times symbol \times if it ends up being the best point for an automatic break. This is analogous to the way the soft hyphen (U+00AD) works in ordinary text.

Interequation alignment can be accomplished by inserting `&`'s in front of the operators, one per equation and not inside math objects, to be aligned at the same horizontal position. For example, the lines

```
a&=b+c
x+y&=3
```

build up as

$$a = b + c$$

$$x + y = 3$$

See also Sec. [3.19](#) on the equation array for similar functionality.

3.24 Size Overrides

UnicodeMath has a command to override the default character sizing. The inverted F character `ƒ` (U+2132) followed by various ASCII characters changes the "font" of the text. For example, `a_ƒA2` builds up as a_2 in contrast to `a_2`, which builds up as a_2 . The

subscript 2 is larger than normal in the former. Some of the \d codes are defined in the table

\dA	One size larger
\dB	Two sizes larger
\dC	One size smaller
\dD	Two sizes smaller

These values are handy for roundtripping increase/decrease argument size context-menu options.

3.25 Math Color

In slide presentations and elsewhere, it can be handy to have math text color and math background color. In fact, Presentation MathML has the [attributes](#) “math-color” and “math-background”. For these, we use the operators Noah Doering added to his [UnicodeMathML](#) implementation to enter text and background.

Inputting math using UnicodeMath into Word and other applications typically doesn’t require special text and background color syntax since these applications have user interfaces (UI) to color text. But UnicodeMath is useful for persisting math in files and in plain-text scenarios that don’t have rich UI. I use UnicodeMath to represent math in many math tests because it’s so much easier to input and read than MathML, LaTeX, RTF, and OMML. So, it’s valuable to be able to include text coloring in UnicodeMath.

Following the UnicodeMath tradition of using non-math UnicodeMath symbols to represent operations, Noah chose \color (U+270E) as the text color operator and \background (U+2601) for background color. The syntax resembles that for n^{th} root $\sqrt[n]{x}$. For example, “ $a + \color{red}q^2 + \background{lime}r^2$ ” builds up to

$$a + q^2 + r^2$$

You can specify a color using the #RRGGBB or #RGB syntax or by the name of a standard HTML color (aqua, black, blue, fuchsia, gray, green, lime, maroon, navy, olive, purple, red, silver, teal, white, yellow).

4. Input Methods

In view of the large number of characters used in mathematics, it is useful to give some discussion of input methods. The ASCII math symbols are easy to find, e.g., + - / * [] () { }, but often need to be used as themselves. To handle these cases and to provide convenient entry of many other symbols, one can use an escape character, the backslash (\), followed by the desired operator or its autocorrect name. Note that a particularly valuable use of UnicodeMath is for inputting formulas into technical documents or programs. In contrast, the direct input of tagged formats like MathML is cumbersome by hand.

4.1 Character Translations

From syntax and typographical points of view, the Unicode minus sign (U+2212) is displayed instead of the ASCII hyphen-minus (U+002D) and the prime (U+2032) is used instead of the ASCII apostrophe (U+0027), but in math zones the minus sign and prime can be entered using these ASCII counterparts. Note that for proper typography, the prime should have a large glyph variant that when superscripted looks correct. The primes in most fonts are chosen to look approximately like a superscript, but they don't provide the desired size and placement to merge well with other superscripts.

Similarly, it is easier to type ASCII letters than italic letters, but mathematical variables are traditionally italicized in print. Accordingly, in math zones translate letters deemed to be variables to the corresponding math alphabetic characters (in the [Mathematical Alphanumeric Symbols Block](#) U+1D400–U+1D7FF or in the [Letterlike Symbols Block](#) U+2100—U+213F) as you can see in [UnicodeMathML](#). Letter combinations corresponding to standard function names like “sin” and “tan” should be represented by ASCII alphabets. As such they are not italicized and are rendered with normal typography, i.e., not mathematical typography. Other post-entry enhancements include mappings like

!!	!!	U+203C
+-	±	U+00B1
-+	∓	U+2213
::	::	U+2237
:=	:=	U+2254
<=	≤	U+2264
>=	≥	U+2265
<<	≪	U+226A
>>	≫	U+226B
~=	≅	U+2245
->	→	U+2192

The pair <- shouldn't map into ←, since expressions like $x < -b$ are common. Also it's not a good idea to map != into ≠, since ! is often used in mathematics to mean factorial, and /= is as easy to type.

Negated counterparts to common mathematical operators can be entered by typing a / in front of the operator by. Operators with this behavior include those in the following table

Operator	Negated op	Input
<	≠	/<
=	≠	/=
>	≠	/>

∃	∄	\backslash exists
∈	∉	\backslash in
⊃	⊄	\backslash ni
~	≈	\backslash sim
≈	≉	\backslash simeq
≅	≇	\backslash cong
≈	≈	\backslash approx
≈	≉	\backslash asymp
≡	≢	\backslash equiv
≤	≰	\backslash le
≥	≱	\backslash ge
≲	≳	\backslash lessgtr
≴	≵	\backslash gtrless
≈	≉	\backslash succeq
⋖	⋗	\backslash prec
⋘	⋙	\backslash succ
≲	≳	\backslash preceq
⊂	⊄	\backslash subset
⊃	⊄	\backslash supset
⊆	⊄	\backslash subseteq
⊇	⊄	\backslash supseteq
⊆	⊄	\backslash sqsubseteq
⊇	⊄	\backslash sqsupseteq

All of these characters are in the U+22xx Unicode block (Mathematical Operators) except for the ASCII characters <, =, and >.

If you don't like an automatic translation when entering math, you can undo the translation by typing, for example, Ctrl+z. Suffice it to say that intelligent input algorithms can dramatically simplify the entry of mathematical symbols and expressions.

4.2 Math Keyboards

A special math shift facility for keyboard entry could bring up proper math symbols. The values chosen can be displayed on an on-screen keyboard. For example, the left Alt key could access the most common mathematical characters and Greek letters, the right Alt key could access italic characters plus a variety of arrows, and the right Ctrl key could access script characters and other mathematical symbols. The numeric keypad offers locations for a variety of symbols, such as sub/superscript digits using the left Alt key. Left Alt CapsLock could lock into the left-Alt symbol set, etc. This approach yields what one might call a "sticky" shift. Other possibilities involve the Num-Lock and ScrollLock keys in combinations with the left/right Ctrl/Alt keys. Pretty soon

one realizes that this approach rapidly approaches literally billions of combinations, that is, several orders of magnitude more than Unicode can handle!

4.3 Hexadecimal Input

A handy hex-to-Unicode entry method can be used to insert Unicode characters in general and math characters in particular. Basically, one types a character's hexadecimal code (in ASCII), making corrections as need be, and then types Alt+x. The hexadecimal code is replaced by the corresponding Unicode character. The Alt+x is a toggle, that is, type it once to convert a hex code to a character and type it again to convert the character back to a hex code. Toggling back to the hex code is very useful for figuring out what a character is if the glyph itself doesn't make it clear or for looking up the character properties in the Unicode Standard. If the hex code is preceded by one or more hexadecimal digits, select the desired code so that the preceding hexadecimal characters aren't included in the code. The code can range up to the value 0x10FFFF, which is the highest character in the 17 planes of Unicode. This input method works in Microsoft Word, Outlook, OneNote, and Notepad. All good editors should support it.

4.4 Pull-Down Menus, Ribbons, Context Menus

Pull-down menus and ribbons are popular methods for handling large character sets, but they tend to be slower than keyboard approaches if you know the right keys to type. A related approach is the symbol gallery, often used for emoji, which is an array of symbols either chosen by the user or displaying the characters in a font. Multiple tabs can organize the symbol selections according to subject matter. On-screen keyboards with symbol galleries are valuable for entry of mathematical expressions and of Unicode text in general. Context menus (right-mouse menus) are quite useful since they provide easy access to context-sensitive options, such as converting a stacked fraction into a linear fraction.

4.5 Macros

The autocorrect and keyboard macro features of some word processing systems provide other ways of entering mathematical characters for people familiar with TeX. For example, typing `\alpha` inserts α if the appropriate autocorrect entry is present. This approach is noticeably faster than using menus and is particularly attractive to those with some familiarity with TeX. It can be sped up dramatically via [math auto-complete](#) as you can see by entering TeX control words in [UnicodeMathML](#).

4.6 UnicodeMath Autocorrect List

The UnicodeMath autocorrect list includes most of those defined in Appendix F of *The TeXbook*, like `\alpha` for α , plus a number of others useful for inputting Unicode Math. AsciiMath has a subset of such keywords but omits the leading backslash. Users can [define](#) their own keywords for convenience or preference, such as `\a` for α , which

requires less typing than the official TeX control word `\alpha`. This also allows localization of the control word list. It'd probably be worth adding an option to make the leading backslash optional. That would speed up keyboard entry of UnicodeMath via math autocorrect. The following table shows the default math autocorrect entries for some structures

Control word	Character	Control word	Character
<code>\int</code>	∫ (U+222B)	<code>\oint</code>	∮ (U+222E)
<code>\sum</code>	∑ (U+2211)	<code>\prod</code>	∏ (U+220F)
<code>\funcapply</code>	(U+2061)	<code>\naryand, \of</code>	⋈ (U+2592)
<code>\rect</code>	◻ (U+25AD)	<code>\sqrt</code>	√ (U+221A)
<code>\open</code>	⊢ (U+251C)	<code>\close</code>	⊣ (U+2524)
<code>\above</code>	⊕ (U+2534)	<code>\below</code>	⊖ (U+252C)
<code>\underbar</code>	— (U+2581)	<code>\overbar</code>	¯ (U+00AF)
<code>\underbrace</code>	⏟ (U+23DF)	<code>\overbrace</code>	⏞ (U+23DE)
<code>\begin</code>	[(U+3016)	<code>\end</code>] (U+3017)
<code>\phantom</code>	⋄ (U+27E1)	<code>\box</code>	◻ (U+25A1)
<code>\hphantom</code>	↔ (U+2B04)	<code>\vphantom</code>	↕ (U+21F3)
<code>\asmash</code>	↑ (U+2B06)	<code>\dsmash</code>	↓ (U+2B07)
<code>\hsmash</code>	↔ (U+2B0C)	<code>\smash</code>	↑ (U+2B0D)
<code>\matrix</code>	■ (U+25A0)	<code>\eqarray</code>	■ (U+2588)

Appendix B contains a default set of keywords containing both *The TeXbook* keywords and UnicodeMath keywords.

4.7 Handwritten Input

Particularly for touch screens, handwritten input is attractive provided the handwriting recognizer is able to decipher the user's handwriting. For this approach, it's desirable to bypass UnicodeMath altogether and recognize built-up mathematical expressions directly.

4.8 Math Dictation

You can say “a squared plus b squared equals c squared” faster than you can write $a^2 + b^2 = c^2$ or type it. UnicodeMath is useful for [math dictation](#) since speech text can be translated into UnicodeMath and then built up. UnicodeMath is significantly closer to math speech than other math formats. You can try it out in English by clicking on the Dictate button in [UnicodeMathML](#).

4.9 Braille

The 6-dot [Nemeth braille encoding](#) was created by Abraham Nemeth for mathematical and scientific notation. It's general enough to encode almost all of UnicodeMath. He started working on his encoding in 1946 and it was first published in 1952 by the American Printing House for the Blind. As such it's the first linear format for math. It's a little like UnicodeMath in that spaces play important roles and it's a globalized notation, so localization isn't needed except for embedded natural language. Also both formats strive to make simple things easy and concise at the cost of additional syntax rules. Because only 64 codes (including the space) are used to encode virtually all of math notation plus a variety of other things, the semantics of the codes depend heavily on their contexts. This level of complexity contrasts with UnicodeMath which has the luxury of the exhaustive Unicode math symbol set. Accordingly, encoding math expressions can become quite tricky as revealed in the [full specification](#). For a less daunting intro, see this [Nemeth Code Cheat Sheet](#). Nemeth recounts some history in this 1991 [interview](#). For Nemeth braille sequences corresponding to Unicode symbols, see [Unicode Math Braille Sequences - Math in Office \(microsoft.com\)](#). You can see the braille for the current equation by clicking on the Braille button in [UnicodeMathML](#).

5. Recognizing Mathematical Expressions

UnicodeMath expressions can be used “as is” for simple documentation purposes. Use in more elegant documentation and in programming languages requires knowledge of the underlying mathematical structure. This section describes some heuristics that can distill math structure out of plain text and gives an algorithm for doing so.

Note that if explicit math-zone-on and math-zone-off characters are desired, Sec. [3.20](#) specifies that [(U+2045) starts a math zone and] (U+2046) ends it. These are not ordinarily used in technical documents. If they do need to be included in a math zone, they can be preceded by the “quote” character \ as described in Sec. [3.2](#).

Many mathematical expressions identify themselves as mathematical, obviating the need to declare them explicitly as such. One well-known TeX problem is TeX's inability to detect expressions that are clearly mathematical, but that are not enclosed within \$'s. If one leaves out a \$ by mistake, one gets many error messages because TeX interprets subsequent text in the wrong mode. This problem is alleviated in LaTeX, which has different math zone start and end delimiters.

An advantage of recognizing mathematical expressions without math-on and math-off syntax is that it is more tolerant to user errors of this sort. Resyncing is automatic, while in TeX one basically has to start up again from the omission in question. Furthermore, this approach could be useful in recognizing and converting the mathematical literature that is not yet available in an object-oriented machine-readable form, into that form.

The basic idea is that math characters identify themselves as such *and* potentially identify their surrounding characters as math characters as well. For example, the fraction / (U+2044) and ASCII slashes, symbols in the range U+2200 through U+22FF, the symbol combining marks (U+20D0..U+20FF), the [math alphanumerics](#) (see U+1D400..U+1D7FF, U+2100..U+214F), and in general, Unicode characters with the mathematics property, identify the characters immediately surrounding them as parts of math expressions.

Most single ASCII letters and all math alphanumerics can start a math zone and reside in one. In English the only exceptions are ‘a’, ‘A’, and ‘I’, which *may* start a math zone. These three letters *do* start a math zone if the character that follows is an operator, a comma, or a period. Similarly, a number or a left parenthesis may start a math zone depending on the character that follows. The algorithm relies on a math dictionary that defines operator properties like unary, binary, relational, and opening/closing brace. Note that identifying a math zone doesn’t imply that the math itself is valid; validation is a separate process. But what’s identified as a math zone is unlikely to be natural language. The following subsections explain the math recognition algorithm.

5.1 Characters that may start a math zone

Scan the text looking for one of the following entities which *may* start a math zone depending on the character(s) that follow:

- Single ‘a’, ‘A’, or ‘I’
- Span of digits, perhaps with an embedded ‘.’, i.e., numbers
- Span of English letters
- One of “(+-[{“

When a character may start a math zone, save its character position (cp) in cpMZ as the potential start of a math zone. Bypass a span of digits or a span of letters. Except for a letter span, set a flag fMightBeMathZone = true indicating that a math zone might be starting.

5.2 Characters that do start a math zone

While scanning the text as in Sec. 5.1, also look for the following characters which *do* start a math zone:

- Single upper- and lower-case letters of the English alphabet other than ‘a’, ‘A’, or ‘I’
- Unicode [math alphanumeric characters](#)
- Unary operators other than ‘+’ and ‘-’, e.g., ‘±’ and ‘∓’
- Opening braces except for ‘(’, ‘[’, and ‘{’
- Symbols in “|ħ[√C∂∃#Δ∇∏∏Σ-”
- Integrals
- Greek letters
- U+2A00..U+2A1C
- ‘|’ (U+2045) which is the UnicodeMath math-zone start delimiter

- Letter spans followed by ‘^’ and math function names

Math function names include the trigonometric functions defined by the template

$$[a[rc]]\{\sin, \cos, \tan, \sec, \csc, \cot, \operatorname{ctg}, \operatorname{tg}\}[h]$$

Other math function names include Pr, arg, deg, det, dim, erf, erfc, exp, gcd, hom, inf, ker, lg, lim, liminf, limsup, ln, log, max, min, sgn, and sup.

5.3 Characters that confirm a possible math zone start

If the flag `fMightBeMathZone` is true, any of the following characters confirm that the math zone has started

- Opening brace
- Binary operator
- Relational operator
- Unary operator
- U+2190..U+22FF
- Greek letter
- Math alphanumeric
- Any in “!-^_|\h”
- U+2A00..U+2AFF

5.4 Spaces and vertical bars inside a math zone

Inside math zones, maintain two counters: `cParen` for parentheses and `cDirac` for bra-ket $\langle \rangle$ expressions. For example, ‘(’ increments `cParen` and ‘)’ decrements it. In the following, `cp` is the current character position and `cpMZ` is the character position at the start of the math zone.

If a character is an ASCII space (U+0020), delete it if

- `cp > cpMZ + 1` and the previous character is unary-binary, binary, or relational
- Following character is a ‘||’ (U+2016)
- Following character is unary-binary, binary, relational
- Following character is any opening brace except for ‘(’

If a character is ‘|’ and the previous character is ‘|’, replace the two ‘|’ by a ‘||’ and delete the space (if any) preceding the resulting ‘||’.

5.5 Characters that end a math zone

A math zone is terminated by

- Unrecognized letter span unless preceded by a binary or relational op, or by a number, or followed by ‘^’
- Characters with codes `< U+0020` or ‘?’
- ‘]’ (U+2046) which is the UnicodeMath math-zone end delimiter
- ASCII space if the previous character isn’t a ‘||’ and either isn’t a ‘,’ or if so, `cParen = cDirac = 0`.

Unless the math zone is terminated by ']', don't include trailing spaces, ',', ':', or '(' in the math zone.

5.6 Example

The algorithm above reformats the plain text

The solution to the general quadratic equation of the form $ax^2 + bx + c = 0$ is given by the quadratic formula:

$$x = (-b \pm \sqrt{b^2 - 4ac}) / 2a$$

where a , b , and c are the coefficients of the equation, and x is the variable representing the solutions. The symbol \pm represents that there are two solutions, one with the plus sign and one with the minus sign.

It's important to note that the quadratic formula only applies to quadratic equations, and may not give a solution for certain values of the coefficients, such as when the discriminant $b^2 - 4ac$ is negative. In these cases, the solutions are complex numbers.

with appropriate math zones as

The solution to the general quadratic equation of the form $ax^2 + bx + c = 0$ is given by the quadratic formula:

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

where a , b , and c are the coefficients of the equation, and x is the variable representing the solutions. The symbol \pm represents that there are two solutions, one with the plus sign and one with the minus sign.

It's important to note that the quadratic formula only applies to quadratic equations, and may not give a solution for certain values of the coefficients, such as when the discriminant $b^2 - 4ac$ is negative. In these cases, the solutions are complex numbers.

The equation $x^2 + bx + c = 0$ and the variables a , b , c , and x are rendered in math zones. The solution for x is rendered as a displayed equation.

6. Using UnicodeMath in Programming Languages

In the middle 1950's, the authors of FORTRAN named their computer language after FORMula TRANslation, but they only went part way. Arithmetic expressions in Fortran and other current high-level languages still do not look like mathematical formulas and considerable human coding effort is needed to translate formulas into their machine comprehensible counterparts. For example, Fortran's superscript construct a^{**k} isn't as readable as a^k and Fortran's subscript $a(k)$ isn't as readable as a_k . Bertrand Russell once said⁷ "a good notation has a subtlety and suggestiveness which at times make it seem almost like a live teacher...and a perfect notation would be a substitute for thought." From this point of view, popular modern computer languages are badly lacking. At least Java allows many Unicode characters as variable names.

Using real mathematical expressions in computer programs would be far superior in terms of readability, reduced coding times, program maintenance, and streamlined documentation. In studying computers we have been taught that this ideal is unattainable, and that one must be content with the arithmetic expression as it is or some other non-mathematical notation such as TeX's. It's worth reexamining this premise. Whereas true mathematical notation clearly used to be beyond the capabilities of machine recognition, we're getting closer now.

In general, mathematics has a wide variety of notations, none of which look like the arithmetic expressions of programming languages. Although ultimately it would be desirable to be able to teach computers how to understand all mathematical expressions, we start with UnicodeMath.

6.1 Advantages of UnicodeMath in Programs

In raw form, these expressions look very like traditional mathematical expressions. With use of the heuristics described above, they can be printed or displayed in traditional built-up form. On disk, they can be stored in pure-ASCII program files accepted by standard compilers and symbolic manipulation programs like Maple, Mathematica, and Macsyma. The translation between Unicode symbols and the ASCII names needed by ASCII-based compilers and symbolic manipulation programs can be carried out via table-lookup (on writing to disk) and hashing (on reading from disk) techniques.

Hence formulas can be both printable in manuscripts *and* computable, either numerically or analytically. Note that this is a goal of MathML as well but attained in a relatively complex way using specialized tools. The idea here is that regular programming languages can have expressions containing standard arithmetic operations and special characters, such as Greek, italics, script, and various mathematical symbols like the square root. Two levels of implementation are envisaged: scalar and vector. Scalar operations can be performed on traditional compilers such as those for C and Fortran. The scalar multiply operator is represented by a raised dot, a legitimate mathematical symbol, instead of the asterisk. To keep auxiliary code to a minimum, the vector implementation requires an object-oriented language such as C++.

The advantages of using UnicodeMath are at least threefold:

- 1) Many formulas in document files can be programmed simply by copying them into a program file and inserting appropriate multiplication dots. This dramatically reduces coding time and errors.
- 2) The use of the same notation in programs and the associated journal articles and books is self documenting. In fact, since many programmers document their programs poorly or not at all, this enlightened choice of notation can change nearly useless or nonexistent documentation into excellent documentation.
- 3) In addition to providing useful tools for the present, these proposed initial steps should help us figure out how to teach computers to understand and

use arbitrary mathematical expressions. Such machine comprehension would greatly facilitate future computations as well as the conversion of the existing paper literature and handwritten input into machine usable form.

The concept is portable to any environment that supports Unicode, and it takes advantage of the fact that high-level languages like C and Fortran accept an “escape” character (“_” and “\$”, respectively) that can be used to access extended symbol sets in a fashion similar to TeX. In addition, the built-in C preprocessor allows niceties such as aliasing the asterisk with a raised dot, which is a legitimate mathematical symbol for multiplication. The Java, C++, and C# languages allow direct use of Unicode variable names, which is a major step in the right direction.

6.2 Comparison of Programming Notations

To get an idea as to the differences between the standard way of programming mathematical formulas and the proposed way, compare the following versions of a C++ routine entitled IHBMWM (inhomogeneously broadened multiwave mixing)

```
void IHBMWM(void)
{
    gammap = gamma*sqrt(1 + I2);
    epsilon = cmplx(gamma+gamma1, Delta);
    alphainc = alpha0*(1-(gamma*gamma*I2/gammap)/(gammap + epsilon));

    if (!gamma1 && fabs(Delta*T1) < 0.01)
        alphacoh = -half*alpha0*I2*pow(gamma/gammap, 3);
    else
    {
        Gamma = 1/T1 + gamma1;
        I2sF = (I2/T1)/cmplx(Gamma, Delta);
        betap2 = epsilon*(epsilon + gamma*I2sF);
        beta = sqrt(betap2);
        alphacoh = 0.5*gamma*alpha0*(I2sF*(gamma + epsilon)
            /(gammap*gammap - betap2))
            *((1+gamma/beta)*(beta - epsilon)/(beta + epsilon)
            - (1+gamma/gammap)*(gammap - epsilon)/
            (gammap + epsilon));
    }
    alpha1 = alphainc + alphacoh;
}
```

```

void IHBMWM(void)
{
     $\gamma = \gamma \cdot \sqrt{1 + I_2};$ 
     $v = \gamma + \gamma_1 + i \cdot \Delta;$ 
     $\alpha_{\text{inc}} = \alpha_0 \cdot (1 - (\gamma \cdot \gamma \cdot I_2 / \gamma') / (\gamma' + v));$ 
    if (! $\gamma_1$  || fabs( $\Delta \cdot T_1$ ) < 0.01)
         $\alpha_{\text{coh}} = -.5 \cdot \alpha_0 \cdot I_2 \cdot \text{pow}(\gamma / \gamma', 3);$ 
    else
    {
         $\Gamma = 1/T_1 + \gamma_1;$ 
         $I_2\mathcal{F} = (I_2/T_1) / (\Gamma + i \cdot \Delta);$ 
         $\beta^2 = v \cdot (v + \gamma \cdot I_2\mathcal{F});$ 
         $\beta = \sqrt{\beta^2};$ 
         $\alpha_{\text{coh}} = .5 \cdot \gamma \cdot \alpha_0 \cdot (I_2\mathcal{F}(\gamma + v) / (\gamma' \cdot \gamma' - \beta^2))$ 
             $\times ((1 + \gamma/\beta) \cdot (\beta - v) / (\beta + v) - (1 + \gamma/\gamma') \cdot (\gamma' - v) / (\gamma' + v));$ 
    }
     $\alpha_1 = \alpha_{\text{inc}} + \alpha_{\text{coh}};$ 
}

```

The above function runs fine with C++ compilers, but C++ does impose some serious restrictions based on its limited operator table. For example, vectors can be multiplied together using dot, cross, and outer products, but there's only one asterisk to overload in C++. In built-up form, the function looks even more like mathematics, namely

```

void IHBMWM(void)
{
     $\gamma = \gamma \cdot \sqrt{1 + I_2};$ 
     $v = \gamma + \gamma_1 + i \cdot \Delta;$ 
     $\alpha_{\text{inc}} = \alpha_0 \cdot \left(1 - \frac{\gamma \cdot \gamma \cdot I_2 / \gamma'}{\gamma' + v}\right);$ 
    if (! $\gamma_1$  || fabs( $\Delta \cdot T_1$ ) < 0.01)
         $\alpha_{\text{coh}} = -.5 \cdot \alpha_0 \cdot I_2 \cdot (\gamma / \gamma')^3;$ 
    else
    {
         $\Gamma = 1/T_1 + \gamma_1;$ 
         $I_2\mathcal{F} = \frac{I_2/T_1}{\Gamma + i \cdot \Delta};$ 
         $\beta^2 = v \cdot (v + \gamma \cdot I_2\mathcal{F});$ 
         $\beta = \sqrt{\beta^2};$ 
         $\alpha_{\text{coh}} = .5 \cdot \gamma \cdot \alpha_0 \cdot \frac{I_2\mathcal{F}(\gamma + v)}{\gamma' \cdot \gamma' - \beta^2} \left( \left(1 + \frac{\gamma}{\beta}\right) \cdot \frac{\beta - v}{\beta + v} - \left(1 + \frac{\gamma}{\gamma'}\right) \cdot \frac{\gamma' - v}{\gamma' + v} \right);$ 
    }
}

```

```

    }
     $\alpha_1 = \alpha_{\text{inc}} + \alpha_{\text{coh}};$ 
}

```

The ability to use the second and third versions of the function was built into the PS Technical Word Processor⁸ circa 1988. With it we already came much closer to true formula translation on input, and the output is displayed in standard mathematical notation. Lines of code could be previewed in built-up format, complete with fraction bars, square roots, and large parentheses. To code a formula, one copies it from a technical document, pastes it into a program file, inserts appropriate raised dots for multiplication and compiles. No change of variable names is needed. Call that 70% of true formula translation! In this way, the C++ function on the preceding page compiles without modification. The code appears nearly the same as the formulas in print [see Chaps. 5 and 8 of Meystre and Sargent⁹].

Questions remain such as whether subscript expressions in UnicodeMath should be treated as part of program-variable names, or whether they should be translated to subscript expressions in the target programming language. Similarly, it would be straightforward to automatically insert an asterisk (indicating multiplication) between adjacent symbols, rather than have the user do it. However here there is a major difference between mathematics and computation: symbolically, multiplication is infinitely precise and infinitely fast, while numerically, it takes time and is restricted to a binary subset of the rationals with limited (although usually adequate) precision. Consequently, for the moment, at least, it seems wiser to consider adjacent symbols as part of a single variable name, just as adjacent ASCII letters are part of a variable name in current programming languages. Perhaps intelligent algorithms will be developed that decide when multiplication should be performed and insert the asterisks optimally.

6.3 Export to TeX

Export to TeX is like export to programming languages but has a modified set of requirements. With current programs, comments are distilled out with distinct syntax. This same syntax can be used in UnicodeMath, although it is interesting to think about submitting a mathematical document to a preprocessor that can recognize and separate out programs for a compiler. In this connection, compiler comment syntax is not particularly pretty; ruled boxes around comments and vertical dividing lines between code and comments are noticeably more readable. So, some refinement of the ways that comments are handled would be very desirable. For example, it would be nice to have a vertical window-pane facility with synchronous window-pane scrolling and the ability to display C code in the left pane and the corresponding // comments in the right pane. Then if one wants to see the comments, one widens the right pane accordingly. On the other hand, to view lines with many characters of code, the // comments needn't get in the way.

With TeX, the text surrounding the mathematics is part and parcel of the technical document, and TeX needs \$'s to distinguish the two. These can be included in the plain text, but it is somewhat ugly. The algorithm described in Sec. 5 goes a long way in determining what is mathematics and what is natural language. Accordingly, the export method consists of identifying the mathematical expressions and enclosing them in \$'s. The special symbols are translated to and from the standard TeX ASCII names as for the program translations. Alternatively one can use LaTeX's `\[...\]` open/close math zone delimiters. You can convert an equation in UnicodeMath to LaTeX in [UnicodeMathML](#) by clicking on the TeX button.

Export to MathML also requires knowing the start and end of a math zone. The built-up functions can all be represented using MathML elements or combinations of elements. The most glaring omission in Presentation MathML is that there's no "*n*-ary" element: one needs to use one of a variety of other elements like `<msub>` along with the desired *n*-ary operator inside an `<mo>`. In addition, one needs to tag numbers, operators, and identifiers.

7. Conclusions

We have shown how with a few additions to Unicode, mathematical expressions can usually be represented with a readable Unicode nearly plain-text format called UnicodeMath. The text consists of combinations of operators and operands. A simple operand consists of a span of non-operators, a definition that substantially reduces the number of parenthesis-override pairs and thereby increases the readability of the plain text. To simplify the notation, operators have precedence values that control the association of operands with operators unless overruled by parentheses. Heuristics can be applied to Unicode math to recognize what parts of a document are math zones. This allows the Unicode plain text to be used in a variety of ways, including in technical document preparation particularly for input purposes, symbolic manipulation, numerical computation, and literals in computer programs.

A variety of syntax choices can be used for a linear format. The choices made for UnicodeMath favor efficient input of mathematical formulae, sufficient generality to support high-quality mathematical typography, the ability to round trip elegant mathematical text at least in a rich-text environment, and a format that resembles a real mathematical notation. Obviously compromises between these goals had to be made.

The heuristics given for recognizing mathematical expressions work well, but they are not infallible. An effective use of heuristics would be by an autoformatting wizard that delimits what it thinks are math zones with on/off codes or a character-format attribute. The user could then overrule any incorrect choices. Once the math zones are identified unequivocally, export to MathML, compilers, and other consumers of mathematical expressions is straightforward.

For further discussion of UnicodeMath and related topics, see the [Math in Office blog](#).

Acknowledgements

This work has benefitted from discussions with many people, notably PS Technical Word Processor users, Asmus Freytag, Barbara Beeton, Ken Whistler, Donald Knuth, Jennifer Michelstein (Halberstam), Ethan Bernstein, Said Abou-Hallawa, Jason Rajtar, Yi Zhang, Geraldine Wade, Ross Mills, John Hudson, Ron Whitney, Richard Lawrence, Sergey Malkin, Alex Gil, Mikhail Baranovsky, Hon-Wah Chan, José Oglesby, Isao Yamauchi, Yuriko Rosnow, Robert Miller, Joe Roni, Jinsong Yu, Sergey Genkin, Victor Kozyrev, Andrei Burago, Eliyezer Kohen, and Michael Aronson. Earlier related work is listed in Ref. 10.

Appendix A. UnicodeMath Grammar

The following grammar is simplified compared to the model in the text. For the general [parsing expression grammar](#) used in [UnicodeMathML](#), see [unicodemathml-parser.pegjs](#).

<i>char</i>	←	Unicode character
<i>space</i>	←	ASCII space (U+0020)
<i>aASCII</i>	←	ASCII A-Z a-z
<i>nASCII</i>	←	ASCII 0-9
<i>anMath</i>	←	Unicode math alphanumeric (U+1D400 – U+1D7FF with some Letterlike symbols U+2102 – U+2134)
<i>anOther</i>	←	Unicode alphanumeric not including <i>anMath</i> nor <i>nASCII</i>
<i>an</i>	←	<i>anMath</i> <i>anOther</i>
<i>diacritic</i>	←	Unicode combining mark
<i>opArray</i>	←	'&' VT '■'
<i>opClose</i>	←)] } '}'
<i>opCloser</i>	←	<i>opClose</i> "\close"
<i>opDecimal</i>	←	'.' ','
<i>opHbracket</i>	←	Unicode math horizontal bracket
<i>opNary</i>	←	Unicode integrals, summation, product, and other nary ops
<i>opOpen</i>	←	(' '[' '{' '⟨'
<i>opOpener</i>	←	<i>opOpen</i> "\open"
<i>opOver</i>	←	'/' "\atop"
<i>opBuildup</i>	←	'_ '^ '√' '∛' '∜' '□' '/' ' ' <i>opArray</i> <i>opOpen</i> <i>opClose</i> <i>opNary</i> <i>opOver</i> <i>opHbracket</i> <i>opDecimal</i>
<i>other</i>	←	<i>char</i> – { <i>an</i> + <i>nASCII</i> + <i>diacritic</i> + <i>opBuildup</i> + CR}
<i>diacriticbase</i>	←	<i>an</i> <i>nASCII</i> '(exp)'
<i>diacritics</i>	←	<i>diacritic</i> <i>diacritics diacritic</i>
<i>atom</i>	←	<i>an</i> <i>diacriticbase diacritics</i>
<i>atoms</i>	←	<i>atom</i> <i>atoms atom</i>

<i>digits</i>	←	<i>nASCII</i> <i>digits nASCII</i>
<i>number</i>	←	<i>digits</i> <i>digits opDecimal digits</i>
<i>expBracket</i>	←	<i>opOpener exp opCloser</i>
	←	' ' <i>exp</i> ' '
	←	' ' <i>exp</i> ' '
<i>word</i>	←	<i>αASCII</i> <i>word αASCII</i>
<i>scriptbase</i>	←	<i>word</i> <i>word nASCII</i> <i>anMath</i> <i>number</i> <i>other</i> <i>expBracket</i> <i>opNary</i>
<i>soperand</i>	←	<i>operand</i> '∞' '-' <i>operand</i> "-∞"
<i>expSubsup</i>	←	<i>scriptbase</i> '_' <i>soperand</i> '^' <i>soperand</i> <i>scriptbase</i> '^' <i>soperand</i> '_' <i>soperand</i>
<i>expSubscript</i>	←	<i>scriptbase</i> '_' <i>soperand</i>
<i>expSuperscript</i>	←	<i>scriptbase</i> '^' <i>soperand</i>
<i>expScript</i>	←	<i>expSubsup</i> <i>expSubscript</i> <i>expSuperscript</i>
<i>entity</i>	←	<i>atoms</i> <i>expBracket</i> <i>number</i>
<i>factor</i>	←	<i>entity</i> <i>entity</i> '!' <i>entity</i> "!!" <i>function</i> <i>expScript</i>
<i>operand</i>	←	<i>factor</i> <i>operand factor</i>
<i>box</i>	←	'□' <i>operand</i>
<i>hbrack</i>	←	<i>opHbracket operand</i>
<i>sqrt</i>	←	'√' <i>operand</i>
<i>cubert</i>	←	'∛' <i>operand</i>
<i>fourthrt</i>	←	'∜' <i>operand</i>
<i>nthrt</i>	←	"√ (" <i>operand</i> '&' <i>operand</i> ")"
<i>function</i>	←	<i>sqrt</i> <i>cubert</i> <i>fourthrt</i> <i>nthrt</i> <i>box</i> <i>hbrack</i>
<i>numerator</i>	←	<i>operand</i> <i>fraction</i>
<i>fraction</i>	←	<i>numerator opOver operand</i>
<i>row</i>	←	<i>exp</i> <i>row</i> '&' <i>exp</i>
<i>rows</i>	←	<i>row</i> <i>rows</i> '@' <i>row</i>
<i>array</i>	←	"\array(" <i>rows</i> ")"
<i>element</i>	←	<i>fraction</i> <i>operand</i> <i>array</i>
<i>exp</i>	←	<i>element</i> <i>exp other element</i>

Appendix B. Character Keywords and Properties

The following table gives the default math keywords, their target characters and codes along with spacing and linear-format build-up properties. A full keyword consists of a backslash followed by a keyword in the table. Some of the keywords may not be in the standard Microsoft 365 math autocorrect list but you can [add](#) them.

Keyword	Glyph	Code	Spacing	UM Property
<code>\above</code>	$\overset{+}{}$	U+2534	ordinary	subsup upper
<code>\abs</code>	(a)	U+249C	ordinary	enclosure
<code>\acute</code>	´	U+0301	ordinary	accent
<code>\aleph</code>	\aleph	U+2135	ordinary	operand
<code>\alpha</code>	α	U+03B1	ordinary	operand
<code>\amalg</code>	\amalg	U+2210	ordinary	nary
<code>\angle</code>	\angle	U+2220	unary	normal
<code>\angmsd</code>	\sphericalangle	U+2221	relational	normal
<code>\angrtvb</code>	\sphericalangle	U+22BE	relational	normal
<code>\angsph</code>	\sphericalangle	U+2222	relational	normal
<code>\aoint</code>	\oint	U+2233	ordinary	nary
<code>\approx</code>	\approx	U+2248	relational	normal
<code>\approxeq</code>	\approx	U+224A	relational	normal
<code>\asmash</code>	\uparrow	U+2B06	ordinary	encl phantom
<code>\ast</code>	*	U+2217	binary	normal
<code>\asymp</code>	\asymp	U+224D	relational	normal
<code>\atop</code>	\atop	U+00A6	ordinary	divide
<code>\backcolor</code>	\backcolor	U+2601	ordinary	enclosure
<code>\backsim</code>	\backsim	U+223D	relational	normal
<code>\backsimeq</code>	\backsimeq	U+22CD	relational	normal
<code>\Bar</code>	$\bar{=}$	U+033F	ordinary	accent
<code>\bar</code>	$\bar{-}$	U+0305	ordinary	accent
<code>\bcancel</code>	$\cancel{}$	U+2572	relational	normal
<code>\because</code>	\because	U+2235	relational	normal
<code>\begin</code>	\lceil	U+3016	open	open
<code>\below</code>	$\underset{-}{}$	U+252C	ordinary	subsup lower
<code>\beta</code>	β	U+03B2	ordinary	operand
<code>\beth</code>	\beth	U+2136	ordinary	operand
<code>\between</code>	\between	U+226C	relational	normal
<code>\bot</code>	\perp	U+22A5	relational	normal
<code>\bigcap</code>	\bigcap	U+22C2	ordinary	nary

Unicode Nearly Plain Text Encoding of Mathematics

<code>\bigcup</code>	U	U+22C2	ordinary	nary
<code>\bigodot</code>	⊙	U+2A00	ordinary	nary
<code>\bigoplus</code>	⊕	U+2A01	ordinary	nary
<code>\bigotimes</code>	⊗	U+2A02	ordinary	nary
<code>\bigsqcap</code>	⊓	U+2A05	ordinary	nary
<code>\bigsqcup</code>	⊔	U+2A06	ordinary	nary
<code>\bigudot</code>	⊙	U+2A00	ordinary	nary
<code>\biguplus</code>	⊕	U+2A04	ordinary	nary
<code>\bigvee</code>	∨	U+22C1	ordinary	nary
<code>\bigwedge</code>	∧	U+22C0	ordinary	nary
<code>\bmatrix</code>	Ⓢ	U+24E2	ordinary	encl matrix
<code>\bowtie</code>	⋈	U+22C8	relational	normal
<code>\bot</code>	⊥	U+22A5	relational	normal
<code>\matrix</code>	■	U+25A0	ordinary	encl matrix
<code>\box</code>	□	U+25A1	ordinary	encl box
<code>\boxdot</code>	◻	U+22A1	binary	normal
<code>\boxed</code>	□	U+25AD	ordinary	encl rect
<code>\boxit</code>	□	U+25AD	ordinary	encl rect
<code>\boxminus</code>	⊖	U+229F	binary	normal
<code>\boxplus</code>	⊕	U+229E	binary	normal
<code>\boxtimes</code>	⊗	U+22A0	binary	normal
<code>\bra</code>	⟨	U+27E8	open	open
<code>\breve</code>	˘	U+0306	ordinary	accent
<code>\bullet</code>	•	U+2219	binary	normal
<code>\bumpeq</code>	≈	U+224F	relational	normal
<code>\cancel</code>	/	U+2571	ordinary	enclosure
<code>\Cap</code>	⋈	U+22D2	binary	normal
<code>\cap</code>	∩	U+2229	binary	normal
<code>\card</code>	©	U+24D2	unary	enclosure
<code>\cases</code>	©	U+24B8	ordinary	encl eqarray
<code>\cbrt</code>	∛	U+221B	open	encl root
<code>\cdot</code>	•	U+22C5	binary	normal

Unicode Nearly Plain Text Encoding of Mathematics

<code>\cdots</code>	…	U+22EF	ordinary	normal
<code>\check</code>	ˇ	U+030C	ordinary	accent
<code>\chi</code>	χ	U+03C7	ordinary	operand
<code>\choose</code>	(c)	U+249E	ordinary	divide
<code>\circ</code>	◦	U+2218	binary	normal
<code>\circeq</code>	≐	U+2257	relational	normal
<code>\circlearrowleft</code>	↺	U+21BA	relational	normal
<code>\circlearrowright</code>	↻	U+21BB	relational	normal
<code>\close</code>	⊥	U+2524	ordinary	close
<code>\clubsuit</code>	♣	U+2663	ordinary	normal
<code>\coint</code>	⌘	U+2232	ordinary	nary
<code>\Colon</code>	::	U+2237	relational	normal
<code>\colon</code>	:	u+2236	relational	normal
<code>\color</code>	✍	u+270E	ordinary	enclosure
<code>\complement</code>	⊂	u+2201	ordinary	normal
<code>\cong</code>	≅	U+2245	relational	normal
<code>\coprod</code>	∏	U+2210	ordinary	nary
<code>\cross</code>	×	U+2A2F	binary	normal
<code>\Cup</code>	⊔	U+22D3	binary	normal
<code>\cup</code>	∪	U+222A	binary	normal
<code>\curlyeqprec</code>	⋈	U+22DE	relational	normal
<code>\curlyeqsucc</code>	⋉	U+22DF	relational	normal
<code>\curlyvee</code>	∨	U+22CE	binary	normal
<code>\curlywedge</code>	∧	U+22CF	binary	normal
<code>\curvearrowleft</code>	↷	U+21B6	relational	normal
<code>\curvearrowright</code>	↶	U+21B7	relational	normal
<code>\cwint</code>	∫	U+2231	ordinary	nary
<code>\dag</code>	†	U+2032	ordinary	unisubsup
<code>\daleth</code>	ℵ	U+2138	ordinary	operand
<code>\dasharrowleft</code>	↔	U+21E0	relational	normal
<code>\dasharrowright</code>	↔	U+21E2	relational	normal
<code>\dashv</code>	⊥	U+22A3	relational	stretch horz

Unicode Nearly Plain Text Encoding of Mathematics

<code>\Dd</code>	\mathcal{D}	U+2145	differential	operand
<code>\dd</code>	\mathcal{d}	U+2146	differential	operand
<code>\ddag</code>	‡	U+2021	ordinary	unisubsup
<code>\ddddot</code>	⋯	U+20DC	ordinary	accent
<code>\dddote</code>	⋯	U+20DB	ordinary	accent
<code>\ddot</code>	¨	U+0308	ordinary	accent
<code>\ddots</code>	⋮	U+22F1	relational	normal
<code>\degc</code>	°C	U+2103	ordinary	operand
<code>\degf</code>	°F	U+2109	ordinary	operand
<code>\degree</code>	°	U+00B0	ordinary	operand
<code>\Delta</code>	Δ	U+0394	ordinary	operand
<code>\Deltaeq</code>	≐	U+225C	ordinary	operand
<code>\delta</code>	δ	U+03B4	ordinary	operand
<code>\det</code>	(v)	U+24B1	ordinary	enclosure matrix
<code>\diamond</code>	◊	U+22C4	binary	normal
<code>\diamondsuit</code>	◇	U+2662	ordinary	normal
<code>\div</code>	÷	U+00F7	binary	normal
<code>\divideontimes</code>	⊗	U+22C7	binary	normal
<code>\dot</code>	·	U+0307	ordinary	accent
<code>\Doteq</code>	≐	U+2251	relational	normal
<code>\doteq</code>	≐	U+2250	relational	normal
<code>\dotminus</code>	⊖	U+2238	binary	normal
<code>\dotplus</code>	⊕	U+2214	binary	normal
<code>\dots</code>	...	U+2026	ordinary	normal
<code>\Downarrow</code>	⇓	U+21D3	relational	normal
<code>\downarrow</code>	↓	U+2193	relational	normal
<code>\downarrows</code>	⇓	U+21CA	relational	normal
<code>\downharpoonleft</code>	↙	U+21C3	relational	normal
<code>\downharpoonright</code>	↘	U+21C2	relational	normal
<code>\dashrightarrow</code>	↞	U+21E2	relational	normal
<code>\dsmash</code>	⤵	U+2B07	ordinary	encl phantom
<code>\ee</code>	ℰ	U+2147	ordinary	operand

Unicode Nearly Plain Text Encoding of Mathematics

<code>\ell</code>	ℓ	U+2113	ordinary	operand
<code>\ellipse</code>	\bigcirc	U+2B2D	ordinary	encl rect
<code>\emptyset</code>	\emptyset	U+2205	unary	operand
<code>\emsp</code>		U+2003	skip	normal
<code>\end</code>	$\}]$	U+3017	close	close
<code>\ensp</code>		U+2002	skip	normal
<code>\epar</code>	$\#$	U+22E5	ordinary	operand
<code>\epsilon</code>	ϵ	U+03F5	ordinary	operand
<code>\eqarray</code>	\blacksquare	U+2588	ordinary	encl eqarray
<code>\eqcirc</code>	\approx	U+2256	relational	normal
<code>\eqgtr</code>	\succ	U+22DD	relational	normal
<code>\eqless</code>	\prec	U+22DC	relational	normal
<code>\eqno</code>	$\#$	U+0023	ordinary	marker
<code>\equiv</code>	\equiv	U+2261	relational	normal
<code>\eta</code>	η	U+03B7	ordinary	operand
<code>\exists</code>	\exists	U+2203	unary	normal
<code>\fallingdotseq</code>	\fallingdotseq	U+2252	relational	normal
<code>\forall</code>	\forall	U+2200	unary	normal
<code>\frown</code>	\frown	U+2322	unary	normal
<code>\funcapply</code>	$\boxed{\circ}$	U+2061	binary	subsup FA
<code>\Gamma</code>	Γ	U+0393	ordinary	operand
<code>\gamma</code>	γ	U+03B3	ordinary	operand
<code>\ge</code>	\geq	U+2265	relational	normal
<code>\geq</code>	\geq	U+2265	relational	normal
<code>\geqq</code>	\cong	U+2267	relational	normal
<code>\gets</code>	\leftarrow	U+2190	ordinary	stretch horiz
<code>\gg</code>	\gg	U+226B	relational	normal
<code>\ggg</code>	\ggg	U+22D9	relational	normal
<code>\gimel</code>	\aleph	U+2137	ordinary	operand
<code>\gneqq</code>	$\not\cong$	U+2269	relational	normal
<code>\gnsim</code>	\approx	U+22E7	relational	normal
<code>\grave</code>	\grave	U+0300	ordinary	accent

Unicode Nearly Plain Text Encoding of Mathematics

<code>\gtrdot</code>	\gtrdot	U+22D7	relational	normal
<code>\gtreqless</code>	\gtreqless	U+22DB	relational	normal
<code>\gtrless</code>	\gtrless	U+2277	relational	normal
<code>\gtrsim</code>	\gtrsim	U+2273	relational	normal
<code>\hairsp</code>		U+200A	skip	normal
<code>\hat</code>	$\hat{}$	U+0302	ordinary	accent
<code>\hbar</code>	\hbar	U+210F	ordinary	operand
<code>\heartsuit</code>	\heartsuit	U+2661	ordinary	normal
<code>\hookleftarrow</code>	\hookleftarrow	U+21A9	relational	stretch horiz
<code>\hookrightarrow</code>	\hookrightarrow	U+21AA	relational	stretch horiz
<code>\hphantom</code>	\hphantom{x}	U+2B04	ordinary	encl phantom
<code>\hsmash</code>	\hsmash{x}	U+2B0C	ordinary	encl phantom
<code>\hvec</code>	\vec{h}	U+20D1	ordinary	accent
<code>\iff</code>	\iff	U+27FF	relational	normal
<code>\ii</code>	ii	U+2148	ordinary	operand
<code>\iiint</code>	\iiint	U+2A0C	ordinary	nary
<code>\iint</code>	\iint	U+222D	ordinary	nary
<code>\iint</code>	\iint	U+222C	ordinary	nary
<code>\Im</code>	\Im	U+2111	ordinary	operand
<code>\imath</code>	\imath	U+0131	ordinary	operand
<code>\in</code>	\in	U+2208	relational	normal
<code>\inc</code>	Δ	U+2206	unary	operand
<code>\infty</code>	∞	U+221E	ordinary	operand
<code>\int</code>	\int	U+222B	ordinary	nary
<code>\intercal</code>	\intercal	U+22BA	binary	normal
<code>\intersection</code>	\cap	U+2229	binary	normal
<code>\iota</code>	ι	U+03B9	ordinary	operand
<code>\itimes</code>	\boxtimes	U+2062	binary	normal
<code>\jj</code>	jj	U+2149	ordinary	operand
<code>\jmath</code>	j	U+0237	ordinary	operand
<code>\kappa</code>	κ	U+03BA	ordinary	operand
<code>\ket</code>	\rangle	U+27E9	close	close

Unicode Nearly Plain Text Encoding of Mathematics

<code>\Lambda</code>	Λ	U+039B	ordinary	operand
<code>\lambda</code>	λ	U+03BB	ordinary	operand
<code>\Langle</code>	⟨⟨	U+27EA	open	open
<code>\langle</code>	⟨	U+27E8	open	open
<code>\lbrack</code>	[[U+27E6	open	open
<code>\lbrace</code>	{	U+007B	open	open
<code>\Lbrack</code>	[[U+27E6	open	open
<code>\lbrack</code>	[U+005B	open	open
<code>\lceil</code>	⌈	U+2308	open	open
<code>\ldiv</code>	/	U+2215	binary	divide
<code>\ldots</code>	...	U+2026	ordinary	normal
<code>\ldsh</code>	↵	U+21B2	ordinary	normal
<code>\le</code>	≤	U+2264	relational	normal
<code>\left</code>	⌊	U+251C	ordinary	normal
<code>\Leftarrow</code>	⇐	U+21D0	relational	stretch horiz
<code>\leftarrow</code>	←	U+2190	relational	stretch horiz
<code>\leftarrowtail</code>	↵	U+21A2	relational	stretch horiz
<code>\leftharpoondown</code>	↷	U+21BD	relational	stretch horiz
<code>\leftharpoonup</code>	↶	U+21BC	relational	stretch horiz
<code>\leftleftarrows</code>	⇐⇐	U+21C7	relational	stretch horiz
<code>\Leftrightarrow</code>	⇔	U+21D4	relational	stretch horiz
<code>\leftrightarrow</code>	↔	U+2194	relational	stretch horiz
<code>\leftrightharpoons</code>	⇌	U+21C6	relational	stretch horiz
<code>\leftrightharpoons</code>	⇌	U+21CB	relational	stretch horiz
<code>\leftrightharpoons</code>	↔↷	U+21AD	relational	stretch horiz
<code>\leftsquigarrow</code>	↵↻	U+21DC	relational	stretch horiz
<code>\leftthreetimes</code>	⋈	U+22CB	binary	normal
<code>\leftwarrow</code>	↵	U+219C	relational	stretch horiz
<code>\leq</code>	≤	U+2264	relational	normal
<code>\leqq</code>	≐	U+2266	relational	normal
<code>\lessdot</code>	⋖	U+22D6	relational	normal
<code>\lesseqgtr</code>	⋗	U+22DA	relational	normal

<code>\lessgtr</code>	\lesssim	U+2276	relational	normal
<code>\lesssim</code>	\lessdot	U+2272	relational	normal
<code>\lfloor</code>	\lfloor	U+230A	open	open
<code>\lhvec</code>	$\vec{\lrcorner}$	U+20D0	ordinary	accent
<code>\ll</code>	\ll	U+226A	relational	normal
<code>\lmoust</code>	\int	U+23B0	relational	normal
<code>\lneqq</code>	\neq	U+2268	relational	normal
<code>\lnot</code>	\neg	U+00AC	relational	normal
<code>\lnsim</code>	\gtrsim	U+22E6	relational	normal
<code>\Longleftarrow</code>	\Lleftarrow	U+27F8	relational	normal
<code>\longleftarrow</code>	\longleftarrow	U+27F5	relational	normal
<code>\Longleftrightarrow</code>	\Lleftrightarrow	U+27FA	relational	normal
<code>\longleftrightharpoonright</code>	\longleftrightarrow	U+27F7	relational	normal
<code>\Longrightarrow</code>	\Rrightarrow	U+27F9	relational	normal
<code>\longrightarrow</code>	\longrightarrow	U+27F6	relational	normal
<code>\looparrowleft</code>	\looparrowleft	U+21AB	relational	normal
<code>\looparrowright</code>	\looparrowright	U+21AC	relational	normal
<code>\lor</code>	\vee	U+2228	binary	normal
<code>\lRhar</code>	\Lrhar	U+21CB	relational	normal
<code>\ltimes</code>	\otimes	U+22C9	binary	normal
<code>\lvec</code>	$\vec{\lrcorner}$	U+20D6	ordinary	accent
<code>\mapsto</code>	\mapsto	U+21A6	relational	stretch horiz
<code>\mapstoleft</code>	\mapstoleft	U+21A4	relational	stretch horiz
<code>\matrix</code>	\blacksquare	U+25A0	ordinary	encl matrix
<code>\medsp</code>		U+205F	ordinary	normal
<code>\mid</code>	\mid	U+2223	relational	list delims
<code>\middle</code>	\circledast	U+24DC	ordinary	normal
<code>\models</code>	\models	U+22A8	relational	stretch horz
<code>\mp</code>	\mp	U+2213	unary/binary	unary/binary
<code>\mu</code>	μ	U+03BC	ordinary	operand
<code>\multimap</code>	\multimap	U+22B8	relational	normal
<code>\nabla</code>	∇	U+2207	unary	operand

Unicode Nearly Plain Text Encoding of Mathematics

<code>\napprox</code>	\approx	U+2249	relational	normal
<code>\naryand</code>	⋈	U+2592	ordinary	normal
<code>\nasymp</code>	\asymp	U+226D	relational	normal
<code>\nbsp</code>		U+00A0	skip	normal
<code>\ncong</code>	\cong	U+2247	relational	normal
<code>\ndiv</code>	\oslash	U+2298	binary	divide
<code>\ne</code>	\neq	U+2260	relational	normal
<code>\nearrow</code>	\nearrow	U+2197	relational	normal
<code>\neg</code>	\neg	U+00AC	unary	normal
<code>\neq</code>	\neq	U+2260	relational	normal
<code>\nequiv</code>	\nequiv	U+2262	relational	normal
<code>\nexists</code>	\nexists	U+2204	unary	normal
<code>\ngeq</code>	\ngeq	U+2271	relational	normal
<code>\ngt</code>	\ngt	U+226F	relational	normal
<code>\ni</code>	\ni	U+220B	relational	normal
<code>\nLeftarrow</code>	\Lleftarrow	U+21CD	relational	stretch horiz
<code>\leftarrow</code>	\leftarrow	U+219A	relational	stretch horiz
<code>\nLeftrightarrow</code>	\Leftrightarrow	U+21CE	relational	stretch horiz
<code>\leftrightarrow</code>	\leftrightarrow	U+21AE	relational	stretch horiz
<code>\nleq</code>	\nleq	U+2270	relational	normal
<code>\nless</code>	\lessdot	U+226E	relational	normal
<code>\nmid</code>	\mid	U+2224	relational	normal
<code>\norm</code>	\parallel	U+2016	ordinary	open/close
<code>\not</code>	$/$	U+2016	binary	divide
<code>\notin</code>	\notin	U+2209	relational	normal
<code>\notni</code>	\nexists	U+220C	relational	normal
<code>\nparallel</code>	\nparallel	U+2226	relational	normal
<code>\nprec</code>	\prec	U+2280	relational	normal
<code>\npreccurlyeq</code>	\preccurlyeq	U+22E0	relational	normal
<code>\nRrightarrow</code>	\Rrightarrow	U+21CF	relational	stretch horiz
<code>\rrightarrow</code>	\rightarrow	U+219B	relational	stretch horiz
<code>\nsim</code>	\sim	U+2241	relational	normal

Unicode Nearly Plain Text Encoding of Mathematics

<code>\nsimeq</code>	\approx	U+2244	relational	normal
<code>\nqssubseteq</code>	$\not\subseteq$	U+22E2	relational	normal
<code>\nqsupseteq</code>	$\not\supseteq$	U+22E3	relational	normal
<code>\nsub</code>	\subsetneq	U+2284	relational	normal
<code>\nsubseteq</code>	\subsetneqq	U+2288	relational	normal
<code>\nsucc</code>	\succneq	U+2281	relational	normal
<code>\nsucccurlyeq</code>	\succneqq	U+22E1	relational	normal
<code>\nsup</code>	\supsetneq	U+2285	relational	normal
<code>\nsupseteq</code>	\supsetneqq	U+2288	relational	normal
<code>\ntriangleleft</code>	\triangleleftneq	U+22EA	relational	normal
<code>\ntriangleleft</code>	\triangleleftneqq	U+22EA	relational	normal
<code>\ntrianglerighteq</code>	\trianglerightneq	U+22EB	relational	normal
<code>\ntrianglerighteq</code>	\trianglerightneqq	U+22ED	relational	normal
<code>\nu</code>	ν	U+03BD	ordinary	operand
<code>\numsp</code>		U+2007	skip	normal
<code>\nVdash</code>	\nVdash	U+22AD	relational	stretch horz
<code>\nvdash</code>	\nvdash	U+22AC	relational	stretch horz
<code>\nwarrow</code>	\nwarrow	U+2196	relational	normal
<code>\oast</code>	\circast	U+229B	binary	normal
<code>\ocirc</code>	\circcirc	U+229A	binary	normal
<code>\odash</code>	\odashdot	U+229D	binary	normal
<code>\odot</code>	\odot	U+2299	binary	normal
<code>\oeq</code>	\oeq	U+229C	binary	normal
<code>\of</code>	\of	U+2592	ordinary	normal
<code>\oiint</code>	\oiint	U+2230	ordinary	nary
<code>\oiint</code>	\oiint	U+222F	ordinary	nary
<code>\oint</code>	\oint	U+222E	ordinary	nary
<code>\Omega</code>	Ω	U+03A9	ordinary	operand
<code>\omega</code>	ω	U+03C9	ordinary	operand
<code>\ominus</code>	\ominus	U+2296	binary	normal
<code>\oo</code>	∞	U+221E	ordinary	operand
<code>\open</code>	\dagger	U+251C	ordinary	open

<code>\oplus</code>	\oplus	U+2295	binary	normal
<code>\oslash</code>	\oslash	U+2298	binary	normal
<code>\otimes</code>	\otimes	U+2297	binary	normal
<code>\over</code>	/	U+002F	binarynsp	divide
<code>\overbar</code>	-	U+00AF	ordinary	encl overbar
<code>\overbrace</code>	$\overbrace{}$	U+23DE	ordinary	stretch over
<code>\overbracket</code>	$\overbracket{}$	U+23B4	ordinary	stretch over
<code>\overline</code>	-	U+00AF	ordinary	encl overbar
<code>\overparen</code>	$\overparen{}$	U+23DC	ordinary	stretch over
<code>\overshell</code>	$\overshell{}$	U+23E0	ordinary	stretch over
<code>\parallel</code>	\parallel	U+2225	relational	normal
<code>\partial</code>	∂	U+2202	unary	operand
<code>\perp</code>	\perp	U+22A5	relational	normal
<code>\phantom</code>		U+27E1	ordinary	encl phantom
<code>\Phi</code>	Φ	U+03A6	ordinary	operand
<code>\phi</code>	ϕ	U+03D5	ordinary	operand
<code>\Pi</code>	Π	U+03A0	ordinary	operand
<code>\pi</code>	π	U+03C0	ordinary	operand
<code>\pitchfork</code>	\pitchfork	U+22D4	binary	normal
<code>\pm</code>	\pm	U+00B1	unary/binary	unary/binary
<code>\pmatrix</code>	\pmatrix{x}	U+25A8	ordinary	encl matrix
<code>\pppprime</code>	''''	U+2057	ordinary	unisubsup
<code>\ppprime</code>	'''	U+2034	ordinary	unisubsup
<code>\pprime</code>	''	U+2033	ordinary	unisubsup
<code>\prime</code>	'	U+2032	ordinary	unisubsup
<code>\prcue</code>	\preccurlyeq	U+227C	relational	normal
<code>\prec</code>	\prec	U+227A	relational	normal
<code>\preccurlyeq</code>	\preccurlyeq	U+227C	relational	normal
<code>\preceq</code>	\preceq	U+2AAF	relational	normal
<code>\precnsim</code>	\precnsim	U+22E8	relational	normal
<code>\precsim</code>	\precsim	U+227E	relational	normal
<code>\prime</code>	'	U+2032	ordinary	Unisubsup

Unicode Nearly Plain Text Encoding of Mathematics

<code>\prod</code>	\prod	U+220F	ordinary	nary
<code>\propto</code>	\propto	U+221D	relational	normal
<code>\Psi</code>	Ψ	U+03A8	ordinary	operand
<code>\psi</code>	ψ	U+03C8	ordinary	operand
<code>\qdrft</code>	$\sqrt[4]{}$	U+221C	open	encl root
<code>\qed</code>	■	U+220E	ordinary	normal
<code>\quad</code>		U+2003	skip	normal
<code>\Rangle</code>	$\rangle\rangle$	U+27EB	close	close
<code>\rangle</code>	\rangle	U+27E9	close	close
<code>\ratio</code>	$:$	U+2236	relational	normal
<code>\rbrack</code>	\rrbracket	U+27E7	close	close
<code>\rbrace</code>	$\}$	U+007D	close	close
<code>\Rbrack</code>	\rrbracket	U+27E7	close	close
<code>\rbrack</code>	$\}$	U+005D	close	close
<code>\rceil</code>	\lceil	U+2309	close	close
<code>\rddots</code>	\ddots	U+22F0	relational	normal
<code>\rdsh</code>	\lrcorner	U+21B3	relational	normal
<code>\Re</code>	\Re	U+211C	ordinary	operand
<code>\rect</code>	\square	U+25AD	ordinary	encl rect
<code>\rfloor</code>	\rfloor	U+230B	close	close
<code>\rho</code>	ρ	U+03C1	ordinary	operand
<code>\rhvec</code>	$\vec{\rho}$	U+20D1	ordinary	accent
<code>\Rightarrow</code>	\Rightarrow	U+21D2	relational	stretch horiz
<code>\right</code>	\dashrightarrow	U+2524	ordinary	normal
<code>\rightangle</code>	\angle	U+2220	unary	normal
<code>\rightarrow</code>	\rightarrow	U+2192	relational	stretch horiz
<code>\rightarrowtail</code>	\rightarrowtail	U+21A3	relational	stretch horiz
<code>\rightharpoondown</code>	\rightharpoondown	U+21C1	relational	stretch horiz
<code>\rightharpoonup</code>	\rightharpoonup	U+21C0	relational	stretch horiz
<code>\rightleftarrows</code>	\rightleftarrows	U+21C4	relational	stretch horiz
<code>\rightleftharpoons</code>	\rightleftharpoons	U+21CC	relational	stretch horiz
<code>\rightrightarrows</code>	\rightrightarrows	U+21C9	relational	stretch horiz

<code>\rightsquigarrow</code>	\rightsquigarrow	U+21DD	relational	stretch horiz
<code>\rightthreetimes</code>	\llcorner	U+22CC	binary	normal
<code>\rightwavearrow</code>	\rightsquigarrow	U+219D	relational	stretch horiz
<code>\risingdotseq</code>	\doteq	U+2253	relational	normal
<code>\rlhar</code>	\Leftrightarrow	U+21CC	relational	normal
<code>\rmoust</code>	\lrcorner	U+23B1	relational	normal
<code>\root</code>	$\sqrt{\quad}$	U+24AD	open	encl root
<code>\rrect</code>	\square	U+25A2	ordinary	encl rnd rect
<code>\rtimes</code>	\rtimes	U+22CA	binary	normal
<code>\sdiv</code>	$/$	U+2044	binarynsp	divide
<code>\sdivide</code>	$/$	U+2044	binarynsp	divide
<code>\searrow</code>	\searrow	U+2198	relational	normal
<code>\setminus</code>	\setminus	U+2216	binary	normal
<code>\sfrac</code>	$/$	U+2044	binarynsp	divide
<code>\Sigma</code>	Σ	U+03A3	ordinary	operand
<code>\sigma</code>	σ	U+03C3	ordinary	operand
<code>\sim</code>	\sim	U+223C	relational	normal
<code>\simeq</code>	\simeq	U+2243	relational	normal
<code>\smash</code>	\updownarrow	U+2B0D	ordinary	encl phantom
<code>\spadesuit</code>	\spadesuit	U+2660	ordinary	normal
<code>\sqcap</code>	\sqcap	U+2293	binary	normal
<code>\sqcup</code>	\sqcup	U+2294	binary	normal
<code>\sqrt</code>	$\sqrt{\quad}$	U+221A	open	encl root
<code>\sqsubset</code>	\sqsubset	U+228F	relational	normal
<code>\sqsubseteq</code>	\sqsubseteq	U+2291	relational	normal
<code>\sqsupset</code>	\sqsupset	U+2290	relational	normal
<code>\sqsupseteq</code>	\sqsupseteq	U+2292	relational	normal
<code>\star</code>	\star	U+22C6	binary	normal
<code>\Subset</code>	\Subset	U+22D0	relational	normal
<code>\subset</code>	\subset	U+2282	relational	normal
<code>\subseteq</code>	\subseteq	U+2286	relational	normal
<code>\subsetneq</code>	\subsetneq	U+228A	relational	normal

Unicode Nearly Plain Text Encoding of Mathematics

<code>\subsub</code>	\Subset	U+2AD5	relational	normal
<code>\subsup</code>	\S	U+2AD3	relational	normal
<code>\succ</code>	\succ	U+227B	relational	normal
<code>\succcurlyeq</code>	\succcurlyeq	U+227D	relational	normal
<code>\succeq</code>	\succeq	U+2AB0	relational	normal
<code>\succnsim</code>	\succnsim	U+22E9	relational	normal
<code>\succsim</code>	\succsim	U+227F	relational	normal
<code>\sum</code>	Σ	U+2211	ordinary	nary
<code>\Supset</code>	\supseteq	U+22D1	relational	normal
<code>\supset</code>	\supset	U+2283	relational	normal
<code>\supseteq</code>	\supseteq	U+2287	relational	normal
<code>\supsetneq</code>	\supsetneq	U+228B	relational	normal
<code>\supsub</code>	\supsetsubset	U+2AD4	relational	normal
<code>\supsup</code>	\supsetsupset	U+2AD6	relational	normal
<code>\surd</code>	\surd	U+221A	open	encl root
<code>\swarrow</code>	\swarrow	U+2199	relational	normal
<code>\tau</code>	τ	U+03C4	ordinary	operand
<code>\therefore</code>	\therefore	U+2234	relational	normal
<code>\Theta</code>	Θ	U+0398	ordinary	operand
<code>\theta</code>	θ	U+03B8	ordinary	operand
<code>\thicksp</code>		U+2005	skip	normal
<code>\thinsp</code>		U+2006	skip	normal
<code>\tilde</code>	\tilde	U+0303	ordinary	accent
<code>\times</code>	\times	U+00D7	binary nsp	normal
<code>\to</code>	\rightarrow	U+2192	relational	stretch horiz
<code>\top</code>	\top	U+22A4	relational	normal
<code>\triangle</code>	\triangle	U+25B3	binary	normal
<code>\trianglelefteq</code>	\trianglelefteq	U+22B4	relational	normal
<code>\trianglerighteq</code>	\trianglerighteq	U+22B5	relational	normal
<code>\tvec</code>	$\vec{\tau}$	U+20E1	ordinary	accent
<code>\twoheadleftarrow</code>	\twoheadleftarrow	U+219E	relational	normal
<code>\twoheadrightarrow</code>	\twoheadrightarrow	U+21A0	relational	normal

<code>\ubar</code>	$_$	U+0332	ordinary	accent
<code>\underbar</code>	$_$	U+2581	ordinary	encl underbar
<code>\underbrace</code>	$\}$	U+23DF	ordinary	stretch under
<code>\underbracket</code>	$\}$	U+23B5	ordinary	stretch under
<code>\underline</code>	$_$	U+2581	ordinary	encl underbar
<code>\underparen</code>	$\)$	U+23DD	ordinary	stretch under
<code>\undershell</code>	$\)$	U+23E1	ordinary	stretch under
<code>\union</code>	\cup	U+222A	binary	normal
<code>\Uparrow</code>	\Uparrow	U+21D1	relational	normal
<code>\uparrow</code>	\uparrow	U+2191	relational	normal
<code>\Updownarrow</code>	\Updownarrow	U+21D5	relational	normal
<code>\updownarrow</code>	\updownarrow	U+2195	relational	normal
<code>\updownarrows</code>	\updownarrows	U+21C5	relational	normal
<code>\upharpoonleft</code>	\upharpoonleft	U+21BF	relational	normal
<code>\upharpoonright</code>	\upharpoonright	U+21BE	relational	normal
<code>\uplus</code>	\uplus	U+228E	binary	normal
<code>\Upsilon</code>	Υ	U+03A5	ordinary	operand
<code>\upsilon</code>	υ	U+03C5	ordinary	operand
<code>\upuparrows</code>	\upuparrows	U+21C8	relational	normal
<code>\varepsilon</code>	ε	U+03B5	ordinary	operand
<code>\varphi</code>	φ	U+03C6	ordinary	operand
<code>\varkappa</code>	\varkappa	U+03F0	ordinary	operand
<code>\varpi</code>	ϖ	U+03D6	ordinary	operand
<code>\varrho</code>	ϱ	U+03F1	ordinary	operand
<code>\varsigma</code>	ς	U+03C2	ordinary	operand
<code>\vartheta</code>	ϑ	U+03D1	ordinary	operand
<code>\vartriangleleft</code>	\vartriangleleft	U+22B2	relational	normal
<code>\vartriangleright</code>	\vartriangleright	U+22B3	relational	normal
<code>\vbar</code>	\vbar	U+2502	ordinary	list delims
<code>\VDash</code>	\Vdash	U+22AB	relational	stretch horz
<code>\Vdash</code>	\Vdash	U+22A9	relational	stretch horz
<code>\vdash</code>	\vdash	U+22A2	relational	stretch horz

<code>\vdots</code>	⋮	U+22EE	relational	normal
<code>\vec</code>	→	U+20D7	ordinary	accent
<code>\vee</code>	∨	U+2228	binary	normal
<code>\Vert</code>	∥	U+2016	ordinary	open/close
<code>\vert</code>		U+007C	ordinary	open/close
<code>\Vmatrix</code>	(n)	U+24A9	ordinary	enclosure matrix
<code>\vmatrix</code>	(v)	U+24B1	ordinary	enclosure matrix
<code>\vphantom</code>	∕	U+21F3	relational	encl phantom
<code>\vthicksp</code>		U+2004	skip	normal
<code>\Vvdash</code>	⊨	U+22AA	relational	stretch horz
<code>\wedge</code>	∧	U+2227	binary	normal
<code>\widehat</code>	^	U+0302	ordinary	accent
<code>\widetilde</code>	~	U+0303	ordinary	accent
<code>\wp</code>	℘	U+2118	ordinary	operand
<code>\wr</code>	ℳ	U+2240	binary	normal
<code>\xcancel</code>	⌂	U+2573	ordinary	enclosure
<code>\Xi</code>	Ξ	U+039E	ordinary	operand
<code>\xi</code>	ξ	U+03BE	ordinary	operand
<code>\zeta</code>	ζ	U+03B6	ordinary	operand
<code>\zwnj</code>		U+200C	ordinary	normal
<code>\zwsp</code>		U+200B	ordinary	normal

Version Differences

The differences between Version 1 and 2 of this paper are largely cosmetic, but there were enough changes in the latter to merit a new number. Version 2 is mostly implemented in Microsoft Word 2007, where it is referred to as the “linear format”. Typing UnicodeMath in Word 2007 or later results in “formula autobuildup”, that is, automatic conversion to the built-up format of an expression when its syntax becomes unambiguous.

In this document features added in Version 3 are implemented in the Microsoft Office applications Word, PowerPoint, Excel, and OneNote (versions 2010 and later). Typically the additions offer convenience over ways needed in Version 2, but no addition is necessary and the Version 2 syntax remains valid in Version 3. The additions were often inspired by [La]TeX. Examples of simplified input are `\choose` for binomial coefficients, `\cases` for alternative definitions, `\pmatrix` for parenthesized matrices,

`\middle` to define a character as a bracket separator, a simpler `prescript` notation, `\root n\of x` notation for n^{th} roots, equation alignment (see Sec. 3.23), size overrides (see Sec. 3.24), and simple negated operator input (see Sec. 4.1). There are also numerous cosmetic changes.

Version 3.1 is mostly a refining of Version 3.0, bringing several topics up to date and using the name UnicodeMath instead of Unicode linear format. Version 3.2 updates Appendix B, polishes the text, and documents math color and two matrix short cuts.

Version 3.3 gives links to [UnicodeMathML](#), a web applet that lets you enter and convert UnicodeMath to MathML and display the results in traditional math notation.

References

1. *The Unicode Standard* <http://www.unicode.org/versions/latest/>
2. Barbara Beeton, Asmus Freytag, Murray Sargent III, Unicode Technical Report #25 “Unicode Support for Mathematics”, <http://www.unicode.org/reports/tr25>
3. Leslie Lamport, *LaTeX: A Document Preparation System, User’s Guide & Reference Manual*, 2nd edition (Addison-Wesley, 1994; ISBN 1-201-52983-1)
4. Donald E. Knuth, *The TeXbook*, (Reading, Massachusetts: Addison-Wesley 1984)
5. *Mathematical Markup Language (MathML)* <http://www.w3.org/Math/>
6. For example, UnicodeMath is used for keyboard entry of mathematical expressions in Microsoft Word, PowerPoint, OneNote and Excel.
7. Bertrand Russell, in his Introduction to *Tractatus Logico-Philosophicus* by Ludwig Wittgenstein, Routledge and Kegan Paul, London 1922 (also currently available at <http://www.kfs.org/~jonathan/witt/tlph.html>).
8. PS Technical Word Processor, Scroll Systems, Inc. (1989). This word processor used a pre-Unicode version of UnicodeMath.
9. P. Meystre and M. Sargent III (1991), *Elements of Quantum Optics*, Springer-Verlag
10. Some of these ideas were discussed in the following presentations: M. Sargent III, *Unicode, Rich Text, and Mathematics*, 7th International Unicode Conference, San Jose, California, Sept (1995); Murray Sargent III and Angel L. Diaz, *MathML and Unicode*, 15th International Unicode Conference, San Jose, California, Sept (1999); Murray Sargent III, *Unicode Plain Text Encoding of Mathematics*, 16th International Unicode Conference, Amsterdam, Holland, March (2000); Murray Sargent III, *Unicode Support for Mathematics*, 17th International Unicode Conference, San Jose, California, Sept (2000); Murray Sargent III, *Unicode Support for Mathematics*, 22nd International Unicode Conference, San Jose, California, Sept (2002); Murray Sargent III, *Unicode Nearly Plain-Text Encoding of Mathematics*, 26th Internationalization and Unicode Conference, San Jose, California, Sept

(2004). Murray Sargent III, *Editing and Display of Mathematics using Unicode*, 29th Internationalization and Unicode Conference, San Francisco, California, March (2006). Murray Sargent III, *Mathematical Input Methods*, 31st Internationalization and Unicode Conference, San Jose, California, Oct (2007). Murray Sargent III, *Math Editing and Display in Microsoft Office*, 33rd Internationalization and Unicode Conference, San José, California, Sept (2009).

This document was prepared using Microsoft Word with Cambria and Cambria Math fonts.

© 2024 Murray Sargent III and Unicode, Inc. This publication is protected by copyright, and permission must be obtained from Unicode, Inc. prior to any reproduction, modification, or other use not permitted by the [Terms of Use](#). Use of all Unicode Products, including this publication, is governed by the Unicode [Terms of Use](#). The authors, contributors, and publishers have taken care in the preparation of this publication but make no express or implied representation or warranty of any kind and assume no responsibility or liability for errors or omissions or for consequential or incidental damages that may arise therefrom. This publication is provided “AS-IS” without charge as a convenience to users. Unicode and the Unicode Logo are registered trademarks of Unicode, Inc. in the United States and other countries.