UnicodeMath A Nearly Plain-Text Encoding of Mathematics Version 3.1

Murray Sargent III Microsoft Corporation 16-Nov-16

1. Inti	oduction	2
2. Enc	oding Simple Math Expressions	3
2.1	Fractions	4
2.2	Subscripts and Superscripts	6
2.3	Use of the Blank (Space) Character	8
3. Enc	coding Other Math Expressions	8
3.1	Delimiters	8
3.2	Literal Operators	11
3.3	Prescripts and Above/Below Scripts	11
3.4	n-ary Operators	12
3.5	Mathematical Functions	13
3.6	Square Roots and Radicals	14
3.7	Enclosures	14
3.8	Stretchy Characters	15
3.9	Matrices	
3.10	Accent Operators	17
3.11	Differential, Exponential, and Imaginary Symbols	18
3.12	Unicode Subscripts and Superscripts	18
3.13	Concatenation Operators	18
3.14	Comma, Period, and Colon	18
3.15	Ordinary Text Inside Math Zones	19
3.16	Space Characters	
3.17	Phantoms and Smashes	21
3.18	Arbitrary Groupings	
3.19	Equation Arrays	
3.20	Math Zones	
3.21	Equation Numbers	
3.22	UnicodeMath Characters and Operands	
3.23	Equation Breaking and Alignment	
3.24	Size Overrides	
-	ut Methods	
4.1	Character Translations	
4.2	Math Keyboards	
4.3	Hexadecimal Input	
4.4	Pull-Down Menus, Ribbons, Context Menus	
4.5	Macros	
4.6	UnicodeMath Autocorrect List	30

4.7	Handwritten Input	. 31
4.8	Speech Input	
4.9	Braille	. 31
5. Rec	ognizing Mathematical Expressions	
6. Usin	ng UnicodeMath in Programming Languages	. 33
6.1	Advantages of UnicodeMath in Programs	. 33
6.2	Comparison of Programming Notations	. 34
6.3	Export to TeX	. 37
7. Con	clusions	. 37
	ledgements	
	ix A. UnicodeMath Grammar	
Append	ix B. Character Keywords and Properties	. 40
Version	Differences	. 49
	ces	

1. Introduction

With a few conventions, <u>Unicode</u> can encode most mathematical expressions in a readable nearly plain text called *UnicodeMath*. The format is linear, but it can be converted to a built-up format that Microsoft Office applications like Word refer to as "Professional". UnicodeMath is more compact and easier to read than [La]TeX,^{3,4} or MathML.⁵ Unlike those formats, it delegates some rich-text properties like text and background colors, font size, footnotes, comments, hyperlinks, etc., to a higher layer. Although one could extend the notation to include such properties, readability would be reduced. Hence in a rich-text environment, UnicodeMath faithfully represents rich mathematical text, while 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 and function-apply arguments (see Secs. <u>3.4</u> and <u>3.5</u>). These aid in displaying built-up functions with proper typography and they also help to interoperate with math-oriented programs and math speech.

A variety of syntax choices can be used for a linear format. The choices made for UnicodeMath favor a number of 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 real mathematical notation.

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 format, UnicodeMath is the most concise. It represents the simple fraction, one half, by the 3 characters "1/2", whereas typical MathML takes 62 characters (consisting of the <mml:mfrac> entity). 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 Office math ribbon. In Word, the structures are defined in OMML (Office MathML) and built up by Word, while for the other apps, the structures are defined in UnicodeMath and built up by RichEdit. The latter are much faster and the equation data much smaller. A dramatic example is the stacked fraction template (empty numerator over empty denominator). In UnicodeMath, this is given by the single character '/'. In OMML, it's 109 characters! LaTeX is considerably shorter at 9 characters "\frac{}{", but is still 9 times longer than UnicodeMath. <u>AsciiMath</u> 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 Unicode-Math 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 <u>Section 3.20</u>), it can recover the built-up mathematics by passing those zones through the RichEdit UnicodeMath <u>MathBuildUp</u> function. In fact, you can roundtrip RichEdit documents containing math zones through the plain-text editor Notepad and the math zones are preserved.

For interchange of math expressions between arbitrary math-aware programs, MathML and other higher-level languages are preferred. At the present time, conversion between UnicodeMath and other math formats is only implemented in Microsoft applications, although UnicodeMath isn't proprietary.

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. <u>Section 3</u> summarizes the usage of the other constructs along with their relative precedences, which are used to simplify the notation. <u>Section 4</u> discusses input methods. <u>Section 5</u> gives ways to recognize mathematical expressions embedded in ordinary text. <u>Section 6</u> explains how Unicode plain text can be helpful in programming languages. <u>Section 7</u> gives conclusions. The appendices present a simplified <u>UnicodeMath grammar</u> and a partial list of <u>operators</u>.

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 Office applications.

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 <u>The Unicode Standard</u>,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

d

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

For more complicated operands (such as those that include 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 really 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

Three built-up fraction variations are available: the "fraction slash" U+2044 (which one might input by typing \sdiv) builds up to a skewed fraction, the "division slash" U+2215 (\ldiv) builds up to a potentially large linear fraction, and the circled slash \oslash (U+2298, \ndiv) builds up a small numeric fraction (although characters other than digits can be used as well). Three kinds of built-up fractions are illustrated by

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

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 'l' (\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 (see <u>Sec. 3.4</u> for a discussion of the *n*-aryand "glue" operator

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

where $(n \mid k)$ 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 Version 3, 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\choose k", assuming that you type a space after the k. This shortcut is included primarily for compatibility with TeX, since (n¦k) is pretty easy to type.

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.

These 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

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

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 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 $a^{h}b$ 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 and the ^ as a small superscript up arrow, 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 line, 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 $a^{b_c}c$ displays as a_c^b (as does $a_c^{b_c}b$). 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 $W^{3\beta}_{\delta_1\rho_1\sigma_2}$, which can be written in UnicodeMath as $W^3\beta_-\delta_1\rho_1\sigma_2$, where Unicode numeric subscripts are used. In TeX, one types

\$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 $\$\alpha_2^3 \vee \beta_2^3 + \alpha_2^3\$$.

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

$$\begin{split} W_{\delta_1\rho_1\sigma_2}^{*}3\beta = U_{\delta_1\rho_1}^{*}3\beta + 1/8\pi^2 & \int \alpha_1^{*}\alpha_2 \\ U_{\rho_1\sigma_2}^{*}1\beta)/U_{\rho_1\sigma_2}^{*}0\beta \end{split} \\ d\alpha'_2 \left[(U_{\delta_1\rho_1}^{*}2\beta - \alpha'_2 + U_{\delta_1\rho_1}^{*}\beta + 1/8\pi^2 + U_{\delta_1\rho_1}^$$

while the standard TeX version reads as

\$\$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[{U_{\delta_1\rho_1}^{2\beta} - \alpha_2' U_{\rho_1\sigma_2}^{1\beta} \over U_{\rho_1\sigma_2}^{0\beta}} \right] \$\$.

Unicode Technical Note 28

2.3 Use of the Blank (Space) Character

The ASCII space character U+0020 is rarely needed for explicit spacing of builtup text since the spacing around operators should be 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 \alphalpha. 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 for build up. Some other operator like + builds up the whole expression, since the operands are unambiguously terminated by such operators.

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 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 do result in spacing. Additional spacing characters are discussed in Sec. <u>3.16</u>.

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 \vdash and \dashv , respectively. Box drawings characters are used for the open/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

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

which has the UnicodeMath " $|x| = \{ w (\&x" \text{ if } "x \ge 0@-\&x" \text{ if } "x < 0) \exists " (see Sec. 3.19 \text{ for a discussion of the equation-array operator }).$

Because the cases construct is fairly common, TeX has the \cases control word for it. This is implemented in UnicodeMath Version 3 with the \cases operator \mathbb{O} . With this the equation above can be written as " $|x| = \mathbb{O}(\&x")$ if $x \ge 0@-\&x"$ if x < 0", which is still a little strange, but you don't have to type the opening curly brace or the \close character (\dashv).

The open/close delimiters can be used to overrule the normal open/close character of delimiters as in the admittedly strange, but nevertheless sometimes used, expression "]a + b[", which has the UnicodeMath " \vdash]a+b+ [". 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 \vdash needs to be treated as an empty open delimiter when it appears before a delimiter like | or], follow the \vdash by a space to force the open-delimiter interpretation.

To suppress automatic sizing and to choose specific sizes, \vdash is followed by a digit '0' -'4' with the meanings in the following table

Digit	Meaning
0	Don't grow
1	TeX \big
2	TeX ∖Big
3	TeX \bigg
4	TeX \Bigg

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-toright 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 |(|x|-|y|)|. 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 |(a|b-c|d)|.

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 (\parallel or \norm) has the same open/close definitions as the absolute value character \mid except that it's always considered to be a delimiter.

Delimiters can also have separators within them. UnicodeMath Version 2 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 the resulting ambiguities are insurmountable in general. One case using U+007C as a separator that can be deciphered is that of 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 \vdash , i.e., ($\vdash|b$). The vertical bar separator grows in size to match the size of the surrounding brackets. In Version 3, 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 '+', which displays the same glyph in UnicodeMath as in built-up form (aside from a possible size reduction), 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.

UnicodeMath operators always 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, while default control words typically look 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 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 delightful that the operator characters look like the operators they represent, while 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 $(_c^{h})a$ creates the prescripted variable $^{b}_{c}a$. Variables can have both prescripts and postscripts (ordinary subscripts and superscripts).

In UnicodeMath Version 3, you can use a prescript notation similar to 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 the ${}^{b}_{c}a$ example, 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 ${}_{1}F_{1}$. In Version 3, 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 \below (\neg) and \above (\perp), respectively. Hence the expression $\lim_{n\to\infty} a_n$ can be represented by $\lim_{n\to\infty} (n\to\infty)$ a_n. Since the operations det, gcd, inf, lim, lim inf, lim sup, max, min, Pr, and sup are common, their below scripts are also accessible by the

usual subscript operator _. So in display mode, $\lim_{n\to\infty} a_n$ can also be represented by $\lim_{n\to\infty} (n\to\infty) a_n$, which is a little easier to type than $\lim_{n\to\infty} (n\to\infty) a_n$.

Although for illustration purposes, the belowscript examples are shown here inline with the script below, ordinarily this choice is only for display-mode math. When inline, below- and abovescripts entered with _ and ^ are shown as subscripts and superscripts, respectively, as are the limits for *n*-ary operators. When entered with and -, 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. <u>3.8</u> for some examples.

3.4 *n*-ary Operators

n-ary 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. This is done by following the sub/superscripted *n*-ary operator by the naryand concatenation operator \naryand ()) which is U+2592. The operand that follows this operator becomes the *n*-aryand. For example, the expression $\int_0^n a^{-n} x dx/(x^2+a^2)$ has the built up form

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

where $xdx/(x^2+a^2)$ is the integrand and d is the Unicode differential d character U+2146. Unlike with the fraction numerator and denominator, the outermost parentheses of a *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.

To delimit more complicated *n*-aryands without using parentheses or brackets of some kind, use the \begin \end ([] see Sec. <u>3.18</u>) 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. <u>3.5</u>). 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 four of the following, OR'd with any of the next three (which were added in Version 3), along with neither or one of the last two

nLimitsDefault	0
nLimitsUnderOver	1
nLimitsSubSup	2
nUpperLimitAsSuperScript	3
nLimitsOpposite	4
nShowLowLimitPlaceHolder	8
nShowUpLimitPlaceHolder	16
fDontGrowWithContent	64
fGrowWithContent	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 with 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 [] 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, one can use $\sinh\{ucapply(\omega-\omega_0)t$, or enclose the argument in [] delimiters.

Since \funcapply isn't the most intuitive name, \of can be used in function-apply contexts. \of autocorrects to (U+2592-) (U+2592) 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. <u>3.16</u>. If an ordinary ASCII space were used, it would imply build up of 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 $\sqrt{(U+221A, \backslash sqrt)}$, $\sqrt[3]{(U+221B, \backslash cbrt)}$, and $\sqrt[4]{(U+221C, \backslash qdrt)}$. 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 $\sqrt[n]{a+b}$.

In UnicodeMath Version 3, you can obtain $\sqrt[n]{a+b}$ using 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, ${}^{n+1}\sqrt{b+c}$ can be input by \root n+1\of(b+c), which is similar to TeX's \root n+1\of{b+c}.

3.7 Enclosures

To enclose an expression in a rectangle one uses the rectangle operator \Box (U+25AD, \rect) followed by the operand representing the expression. This syntax is similar to that for the square root. For example $\Box (E = mc^2)$ displays as $E = mc^2$. The same approach is used to put an overbar above an expression, namely follow the overbar operator $\overline{(U+00AF, \vee overbar)}$ by the desired operand. For an underbar, use the operator $_$ (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 <menclose> element, except for roots, which are represented as discussed in the previous Section. The general syntax for enclosing an expression x is $\Box(n\&x)$, where nis 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 can also be given by $\Box(2\&x)$ and the underbar by $\Box(8\&x)$.

Other enclosures such as rounded box \Box , circle, long division, actuarial, and ellipse \Box 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 $\Box(n\&x)$, where \Box 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

$$\begin{array}{c}
k \text{ times} \\
\overline{x + \dots + x} \\
\underline{x + y + z} \\
>0
\end{array}$$

The UnicodeMath for these are $(x+\dots+x)^{(k \text{ "times"})}$ and $(x+y+z)_{(>0)}$, 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 belowscript and abovescript operators, respectively. The horizontal stretchable brackets are given in the following table

U+23DC	(\overparen
U+23DD)	\underpa-
	0	ren
U+23DE	}	\overbrace
U+23DF	Ļ	\underbrace
U+23E0	ĺ	\overshell

U+23E1		\undershell
U+23B4	Γ	\over-
		bracket
U+23B5		\under-
		bracket

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

$(a+b)^{\perp} \rightarrow$	$\overrightarrow{a+b}$
$(a+b)_{\top} \rightarrow$	a + b
$\rightarrow -a + b$	$\xrightarrow{a+b}{\longrightarrow}$
$\rightarrow -(a+b)$	$\xrightarrow[a+b]{}$

3.9 Matrices

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

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

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

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 in UnicodeMath Version 3 with the \pmatrix operator (m). So (m)(a&b@c&d) displays as

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

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 it 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., $a' \wedge c$ 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

\hat	U+0302	â
\check	U+030C	ă
\tilde	U+0303	ã
\acute	U+0301	á
\grave	U+0300	à
\dot	U+0307	à
\ddot	U+0308	ä
\dddot	U+20DB	ä
\bar	U+0304	ā
\vec	U+20D7	ā

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, $(a + b)^{\circ}$ renders as 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* (\mathbb{D} , *d*, *e*, *i*, *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 \mathbb{D} , d, e, i, 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 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

All remaining 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. <u>3.16</u> 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

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

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

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. <u>2.3</u>, 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	Autocor-
			rect
0 em	U+200B	zero-width space	\zwsp
1/18 em	U+200A	veryverythinmathspace	\hairsp
2/18 em	U+200A U+200A	verythinmathspace	
3/18 em	U+2009	thinmathspace	\thinsp
4/18 em	U+205F	mediummathspace	\medsp
5/18 em	U+2005	thickmathspace	\thicksp
6/18 em	U+2004	verythickmathspace	\vthicksp
7/18 em	U+2004 U+200A	veryverythickmathspace	
9/18 em	U+2002	ensp	∖ensp
18/18 em	U+2003	emsp	∖emsp
digit width	U+2007	numsp	\numsp
space width	U+00A0	no-break space	\nbsp

In general, spaces act as concatenation operators and cause build up of higherprecedence operators that precede them. But it's useful for the zero-width space (U+200B) to be treated as an operand character and not to cause build up of 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) were used after the "lim", it would imply build up of the "lim" function, rather than being part of the "lim sup" or "lim inf" function.

In math zones, most spacing is automatically 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* and Appendix F of the MathML 2.0 specification)

	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 $4/18^{\text{th}}$ em between the a and the + and between the + and the b. Similarly there is $5/18^{\text{th}}$ 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) is handy for use as a null argument. 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 unkerned 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

Autocor-	LF op	Op name	width	as-	de-	ink
rect				cent	scent	
\phantom	♦ U+27E1	white concave-sided dia- mond	W	а	d	no
\hphantom	⇔ U+2B04	white left-right arrow	W	0	0	no
\vphantom	\$ U+21F3	white up-down arrow	0	а	d	no
\smash	‡ U+2B0D	black up-down arrow	W	0	0	yes
\asmash	1 U+2B06	black up arrow	W	0	d	yes
\dsmash	↓ U+2B07	black down arrow	W	а	0	yes
\hsmash	↔ U+2B0C	black left-right arrow	0	а	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 \hsmash phantom, 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 ets 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

$$10x + 3y = 2$$
$$3x + 13y = 4$$

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

3.20 Math Zones

<u>Section 5</u> discusses heuristic methods 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 [...] and LaTeX inline math zones can have the form (...).

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 $(E=mc^2#(30))$ or more simply just $E=mc^2#(30)$ renders as

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

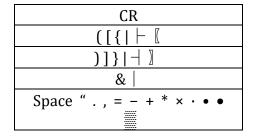
3.22 UnicodeMath Characters and Operands

UnicodeMath divides the roughly 128,000 assigned Unicode characters into three 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, Unicode 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 <u>Appendix</u> <u>A</u>.

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 surprise one at first. Most notably, the space (U+0020) is an important operator in UnicodeMath since it can be used to terminate operands as discussed in Sec. <u>2-3</u>. 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



Unicode Technical Note 28

Unicode Nearly Plain Text Encoding of Mathematics

/
$\int \Sigma \Pi$
_ ^
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

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

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 denomina- tor exp_2 . Numerator and denominator expressions are termi- nated 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. Some- times called a stack.
base^exp1	Superscript expression exp_1 to the base <i>base</i> . The super- scripts $0-9+-()$ exist as Unicode symbols. Sub/superscript ex- pressions are terminated, for example, by /*]) and blank. Sub/superscript operators associate right to left.
base_exp1	Subscript expression exp_1 to the base <i>base</i> . The subscripts $0-9 + -()$ exist as Unicode symbols.
base_exp1^exp2	Subscript expression exp_1 and superscript expression exp_2 to the base <i>base</i> . The subscripts $0-9+-()$ exist as Unicode symbols.

(_exp1^exp2)base	Prescript the subscript exp_1 and superscript exp_2 to the base base.
$base \perp exp_1$	Display expression <i>exp</i> ¹ centered above the base <i>base</i> . Above/below script operators associate right to left.
$base - exp_1$	Display expression <i>exp</i> ¹ centered below the base <i>base</i> .
[<i>exp</i> ₁]	Surround exp_1 with built-up brackets. Similarly for { } and (). Similarly for { }, (), . See Sec. 3.1 for generalizations.
$[exp_1]^{exp_2}$	Surround <i>exp</i> ¹ with built-up brackets followed by super- scripted <i>exp</i> ² (moved up high enough).
$\Box exp_1$	Abstract box around <i>exp</i> ₁ .
$\Box exp_1$	Rectangle around <i>exp</i> ₁ .
<i>_exp</i> 1	Underbar under <i>exp</i> 1 (underbar operator is U+2581, not the ASCII underline character U+005F).
exp_1	Overbar above <i>exp</i> ₁ .
$\sqrt{exp_1}$	Square root of <i>exp</i> ₁ .
∛exp1	Cube root of <i>exp</i> ₁ .
∜ <i>exp</i> 1	Fourth root of <i>exp</i> ₁ .
$\sqrt{(exp_1\&exp_2)}$	$exp_1^{\text{th}} \operatorname{root} \operatorname{of} exp_2.$
$\sum exp_1^exp_2$	Summation from <i>exp</i> ¹ to <i>exp</i> ² with summand <i>exp</i> ³ <i>exp</i> ¹ and ^ <i>exp</i> ² are optional.
$\prod_{exp_1} exp_2 exp_3$	Product from <i>exp</i> ¹ to <i>exp</i> ² with multiplicand <i>exp</i> ³ . <i>_exp</i> ¹ and ^ <i>exp</i> 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.
(<i>exp</i> ₁ [& <i>exp</i> ₂] [@	Align <i>exp</i> ¹ over <i>exp</i> ^{<i>n</i>-1} , etc., to build up an array (see Appendix
	A for a more complete syntax).

 exp_{n-1} [& exp_n]...])

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 blank (U+0020). This blank does not show up when the expression is built up. Blanks that don't terminate operands may be used to space formulas in addition to the built-in spacing provided by a math display engine. Blanks are discussed in greater detail in Sec. <u>2-3</u>.

```
Unicode Technical Note
28
```

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 similar to that for arithmetic expressions, it differs occasionally from the latter.

3.23 Equation Breaking and Alignment

UnicodeMath Version 3 has two features aiding equation breaking and alignment in display math zones. A soft (optional) line break is created by the invisible times (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. 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 × if it ends up being the best point for an automatic break. This is analogous to the way the soft hyphen (U+00AD) is used 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. <u>3.19</u> on the equation array for similar functionality.

3.24 Size Overrides

UnicodeMath Version 3 has a command to override the default character sizing. The inverted F character $\exists (U+2132)$ followed by various ASCII characters changes the "font" of the text. For example, a_ \exists A2 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 \exists codes are defined in the table

AЬ	One size larger
ЯF	Two sizes larger
ЗF	One size smaller
DЬ	Two sizes smaller

These values are handy for roundtripping increase/decrease argument size contextmenu options.

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 very cumbersome if attempted 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 when used as mathematical variables, such letters are traditionally italicized in print. Accordingly a user might want to make italic the default alphabet in a math context, reserving the right to overrule this default when necessary. A more elegant approach in math zones is to translate letters deemed to be standalone to the appropriate math alphabetic characters (in the range U+1D400–U+1D7FF or in the Letterlike Block U+2100–U+213F). Letter combinations corresponding to standard function names like "sin" and "tan" should be represented by ASCII alphabetics. 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
-+	Ŧ	U+2213
::	::	U+2237
:=	:=	U+2254
<=	≤	U+2264
>=	≥	U+2265
<<	~	U+226A
>>	\gg	U+226B
~=	Ш	U+2245

-> → U+2192

The pair <- shouldn't map into \leftarrow , since expressions like x < -b are common. Also it's not a good idea to map != into \neq , since ! is often used in mathematics to mean factorial.

In UnicodeMath Version 3, 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	
<	≮	/<	
=	≠	/=	
>	≯	/>	
Э	∄	/\exists	
E	¢	/\in	
Э	⊅	/\ni	
~	4	/\sim	
~	≄	/\simeq	
≅	≇	/\cong	
≈	≉	/\approx	
X	*	/\asymp	
≡	≢	/\equiv	
	4	/\le	
≥	≱	/\ge	
≶	₩	/\lessgtr	
≷	₩	/\gtrless	
≽	≯	/\succeq	
~	⊀	/\prec	
>	*	/\succ	
<	¥	/\preceq	
	⊄	/\subset	
D	⊅	/\supset	
⊆	⊈	/\subseteq	
⊇	⊉	/\supseteq	
⊑	⊈	/\sqsubseteq	
⊒	⊉	/\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

Computers have multilingual capabilities with keyboards for many different languages. It's desirable to add math keyboards as well. 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 NumLock 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 kind of input is supported, for example, in Microsoft Word and in WordPad.

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. Similarly, one can assign a UnicodeMath expression to a control word. For example, typing \integral in a Microsoft Word math zone inserts

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

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 control words but omits the leading backslash. The user can modify such control words in the Office math autocorrect list or add them explicitly, but 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

Control word	Character	Control word	Character
\int	∫ (U+222B)	\oint	∮ (U+222E)
\sum	∑ (U+2211)	\prod	∏ (U+220F)
\funcapply	(U+2061)	\naryand, \of	(U+2592)
\rect	□ (U+25AD)	\sqrt	√ (U+221A)
\open	⊦(U+251C)	\close	+ (U+2524)
\above	⊥ (U+2534)	\below	т (U+252C)
\underbar	_ (U+2581)	\overbar	⁻ (U+00AF)
\underbrace	∽(U+23DF)	\overbrace	←(U+23DE)
\begin	[(U+3016)	\end] (U+3017)
\phantom	\$(U+27E1)	\box	□ (U+25A1)
\hphantom	⇔(U+2B04)	\vphantom	\$(U+21F3)
\asmash	1 (U+2B06)	\dsmash	↓(U+2B07)
\hsmash	↔(U+2B0C)	\smash	\$ (U+2B0D)
\matrix	■ (U+25A0)	\eqarray	(U+2588)

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

Users can define their own control words for convenience or preference, such as α , which requires less typing than the official TeX control word α . This also allows localization of the control word list.

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 Speech Input

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 nevertheless useful for math speech input since speech text can be translated into UnicodeMath and then built up. UnicodeMath is significantly closer to math speech than other math formats.

4.9 Braille

The 6-dot <u>Nemeth braille encoding</u> was created by Abraham Nemeth for mathematical and scientific notation. It's general enough to encode almost all of Unicode-Math. 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 math linear format. 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. But because a mere 64 codes 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 <u>full specification</u>. For a less daunting intro, see this <u>Nemeth Code Cheat Sheet</u>. Nemeth recounts some history in this 1991 <u>interview</u>.

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 of the heuristics that can distill the structure out of plain text.

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 be 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 La-TeX, 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 much 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 machinereadable form, into that form.

It is possible to use a number of heuristics for identifying mathematical expressions and treating them accordingly. These heuristics are not foolproof, but they lead to the most popular choices. Special commands discussed at the end of this section can be used to overrule these choices. Ultimately the approach could be used as an autoformat style wizard that tags expressions with a rich-text math style whose state is revealed to the user by a toolbar button. The user could then override cases that were tagged incorrectly.

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.

If Latin letter mathematical variables are already given in one of the math alphabets, they are considered parts of math expressions. If they are not, one can still have some recognition heuristics as well as the opportunity to italicize appropriate variables. Specifically ASCII letter pairs surrounded by whitespace are often mathematical expressions, and as such should be italicized in print. If a letter pair fails to appear in a list of common English and European two-letter words, it is treated as a mathematical expression and italicized. Many Unicode characters are not mathematical in nature and suggest that their neighbors are not parts of mathematical expressions.

Strings of characters containing no whitespace but containing one or more unambiguous mathematical characters are generally treated as mathematical expressions. Certain two-, three-, and four-letter words inside such expressions should *not* be italicized. These include trigonometric function names like sin and cos, as well as ln, cosh, etc. Words or abbreviations, often used as subscripts (see the program in Sec. 6), also should not be italicized, even when they clearly appear inside mathematical expressions.

Special cases will always be needed, such as in documenting the syntax itself. The literal operator introduced earlier (\backslash) causes the operator that follows it to be

treated as an nonbuildup operator. This allows the printing of characters without modification that by default are considered to be mathematical and thereby subject to a changed display. Similarly, mathematical expressions that the algorithms treat as ordinary text can be sandwiched between math-on and math-off symbols or by an ordinary text attribute if they need to be embedded in the math zone, e.g., in the numerator of a fraction.

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 a lot closer now.

In general, mathematics has a very 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 at once 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 leads to an unprecedented level of self documentation. In fact, since many programmers document their programs poorly or not at all, this enlightened choice of notation can immediately 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 accomplish the ultimate goal of teaching 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 hand written 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 ith a raised dot, which is a legitimate mathematical symbol for multiplication. The Java and C# languages allow direct use of Unicode variable names, which is a major step in the right direction. Compatibility with unenlightened ASCII-only compilers can be done via an ASCII representation of Unicode characters.

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);
          upsilon = cmplx(gamma+gamma1, Delta);
          alphainc = alpha0*(1-(gamma*gamma*I2/gammap)/(gammap + upsilon));
          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 = upsilon*(upsilon + gamma*I2sF);
                    beta = sqrt(betap2);
                    alphacoh = 0.5*gamma*alpha0*(I2sF*(gamma + upsilon)
                                         /(gammap*gammap – betap2))
                                         *((1+gamma/beta)*(beta – upsilon)/(beta + upsilon)
                                         - (1+gamma/gammap)*(gammap – upsilon)/
                                         (gammap + upsilon));
          }
          alpha1 = alphainc + alphacoh;
}
void IHBMWM(void)
{
          \gamma = \gamma \bullet \sqrt{(1 + I_2)};
          \upsilon = \gamma + \gamma_1 + i \bullet \Delta;
          \alpha_{\text{inc}} = \alpha_0 \bullet (1 - (\gamma \bullet \gamma \bullet I_2 / \gamma') / (\gamma' + v));
          if (! \gamma_1 || fabs(\Delta \bullet T_1) < 0.01)
                    \alpha_{\rm coh} = -.5 \bullet \alpha_0 \bullet I_2 \bullet pow(\gamma/\gamma', 3);
          else
          {
                    \Gamma = 1/T_1 + \gamma_1;
                    I_2 \mathcal{F} = (I_2/T_1)/(\Gamma + i \bullet \Delta);
                    \beta^2 = v \bullet (v + \gamma \bullet I_2 \mathcal{F});
                    \beta = \sqrt{\beta^2};
                    \alpha_{-} \text{coh} = .5 \bullet \gamma \bullet \alpha_{0} \bullet (I_{2} \mathcal{F}(\gamma + v) / (\gamma' \bullet \gamma' - \beta^{2}))
                               \times ((1+\gamma/\beta) \bullet (\beta-\upsilon)/(\beta+\upsilon) - (1+\gamma/\gamma') \bullet (\gamma'-\upsilon)/(\gamma'+\upsilon));
          }
          \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

$$\begin{cases} \gamma = \gamma \cdot \sqrt{1 + I_2}; \\ v = \gamma + \gamma_1 + i \cdot \Delta; \\ \alpha_{inc} = \alpha_0 \cdot \left(1 - \frac{\gamma \cdot \gamma \cdot I_2/\gamma'}{\gamma' + v}\right); \\ if (! \gamma_1 || fabs(\Delta \cdot T_1) < 0.01) \\ \alpha_{coh} = -.5 \cdot \alpha_0 \cdot I_2 \cdot (\gamma/\gamma')^3; \\ else \\ \{ \Gamma = 1/T_1 + \gamma_1; \\ I_2 \mathcal{F} = \frac{I_2/T_1}{\Gamma + i \cdot \Delta}; \\ \beta^2 = v \cdot (v + \gamma \cdot I_2 \mathcal{F}); \\ \beta = \sqrt{\beta^2}; \\ \alpha_{coh} = .5 \cdot \gamma \cdot \alpha_0 \cdot \frac{I_2 \mathcal{F}(\gamma + v)}{\gamma' \cdot \gamma' - \beta^2} \left(\left(1 + \frac{\gamma}{\beta}\right) \cdot \frac{\beta - v}{\beta + v} - \left(1 + \frac{\gamma}{\gamma'}\right) \cdot \frac{\gamma' - v}{\gamma' + v} \right); \\ \} \\ \alpha_1 = \alpha_{inc} + \alpha_{coh}; \end{cases}$$

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 similar to 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 heuristics described in Sec. 5 go 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.

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, which we call UnicodeMath. The text consists of combinations of operators and operands. A Unicode Technical Note

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, and numerical computation.

A variety of syntax choices could be used for a linear format. The choices made in this paper 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 the heuristics would be by an autoformatting wizard that delimits what it thinks are math zones with on/off codes or a characterformat 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 <u>Math in Office</u> <u>blog</u> and Chapter 6 in <u>Creating Research and Scientific Documents with Microsoft</u> <u>Word</u>.

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, 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, and Eliyezer Kohen. Earlier related work is listed in Ref. 10.

Appendix A. UnicodeMath Grammar

This grammar is simplified compared to the model in the text.

char	\leftarrow	Unicode character
space	\leftarrow	ASCII space (U+0020)
αASCII	\leftarrow	ASCII A-Z a-z
nASCII	\leftarrow	ASCII 0-9
αnMath	\leftarrow	Unicode math alphanumeric (U+1D400 – U+1D7FF with some
		Letterlike symbols U+2102 – U+2134)
αnOther	\leftarrow	Unicode alphanumeric not including <i>αnMath</i> nor <i>nASCII</i>

αn diacritic opArray opClose opCloser opDecimal opHbracket opNary opOpen opOpener opOpener opOver opBuildup	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	Unicode combining mark '&' VT '■' ')' ']' '}' ')' <i>opClose</i> "\close" '.' ',' Unicode math horizontal bracket Unicode integrals, summation, product, and other nary ops
diacriticbase diacritics atom atoms digits number	$\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow \uparrow \downarrow$	αn nASCII '(' exp ')' diacritic diacritics diacritic αn diacriticbase diacritics atom atoms atom
expBracket	\leftarrow	opOpener exp opCloser ' ' exp ' ' ' ' exp ' '
word scriptbase	\leftarrow	word word nASCII αnMath number other expBracket
expSubscript	\leftarrow \leftarrow	opNary operand '∞' '-' operand "-∞" scriptbase '_' soperand '^' soperand scriptbase '^' soperand scriptbase '_' soperand scriptbase '^' soperand expSubsup expSubscript expSuperscript
entity factor operand box hbrack sqrt	$\leftarrow \leftarrow$	atoms expBracket number entity entity '!' entity "!!" function expScript factor operand factor '□' operand opHbracket operand '√' operand

cubert fourthrt nthrt function numerator fraction	$\begin{array}{c} \downarrow \\ \downarrow \\ \downarrow \\ \downarrow \end{array}$	'∛' operand '∜' operand "√ (" operand '&' operand ')' sqrt cubert fourthrt nthrt box hbrack operand fraction numerator opOver operand
row rows array	\leftarrow	exp row '&' exp row rows '@' row "\array(" rows ')'
element exp		fraction operand array element exp other element

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.

Keyword	Glyph	Code	Spacing	LF Property
\above	1	U+2534	ordinary	subsup upper
\acute	,	U+0301	ordinary	accent
\aleph	х	U+2135	ordinary	operand
\alpha	α	U+03B1	ordinary	operand
\amalg	Ш	U+2210	ordinary	nary
\angle	۷	U+2220	relational	normal
\aoint	∲	U+2233	ordinary	nary
\approx	≈	U+2248	relational	normal
\asmash	1	U+2B06	ordinary	encl phantom
\ast	*	U+2217	binary	normal
\asymp	Х	U+224D	relational	normal
\atop		U+00A6	ordinary	divide
\Bar	=	U+033F	ordinary	accent
\bar	-	U+0305	ordinary	accent
\because	:	U+2235	relational	normal
\begin	ľ	U+3016	open	open
\below	т	U+252C	ordinary	subsup lower
\beta	β	U+03B2	ordinary	operand

) h - c h	5	11.2126		J
\beth	ב .	U+2136	ordinary	operand
\bot		U+22A5	relational	normal
\bigcap	Π	U+22C2	ordinary	nary
\bigcup	U	U+22C2	ordinary	nary
\bigodot	\odot	U+2A00	ordinary	nary
\bigoplus	\oplus	U+2A01	ordinary	nary
\bigotimes	\otimes	U+2A02	ordinary	nary
\bigsqcup	Ш	U+2A06	ordinary	nary
\biguplus	Ĥ	U+2A04	ordinary	nary
\bigvee	V	U+22C1	ordinary	nary
\bigwedge	Λ	U+22C0	ordinary	nary
\bowtie	×	U+22C8	relational	normal
\bot	T	U+22A5	relational	normal
\box		U+25A1	ordinary	encl box
\bra	<	U+27E8	open	open
\breve	v	U+0306	ordinary	accent
\bullet	•	U+2219	binary	normal
\cap	Ω	U+2229	binary	normal
\cbrt	∛	U+221B	open	encl root
\cdot	•	U+22C5	binary	normal
\cdots		U+22EF	ordinary	normal
\check	~	U+030C	ordinary	accent
\chi	χ	U+03C7	ordinary	operand
\circ	o	U+2218	binary	normal
\close	-1	U+2524	ordinary	close
\clubsuit	÷	U+2663	ordinary	normal
\coint	∲	U+2232	ordinary	nary
\cong	≅	U+2245	relational	normal
\cup	U	U+222A	binary	normal
\daleth	7	U+2138	ordinary	operand
\dashv	Ч	U+22A3	relational	stretch horz
\Dd	D	U+2145	differential	operand
\dd	đ	U+2146	differential	operand
\ddddot	••••	U+20DC	ordinary	accent
\dddot	•••	U+20DB	ordinary	accent
	I		-	

\ddot $0+0308$ ordinaryaccent\ddots \ddots $U+22F1$ relationalnormal\degree \circ $U+00B0$ ordinaryoperand\delta Δ $U+0394$ ordinaryoperand\delta δ $U+03B4$ ordinaryoperand\diamond \diamond $U+22C4$ binarynormal\diamondsuit \diamond $U+2662$ ordinarynormal\diamondsuit \diamond $U+2662$ ordinarynormal\dot $U+00F7$ binarynormal\dot $U+00F7$ binarynormal\dot $U+0307$ ordinaryaccent\dots $U+2250$ relationalnormal\dots $U+2266$ ordinarynormal\dots $U+2266$ ordinarynormal\dots $U+2260$ relationalnormal\dots $U+2260$ ordinarynormal\dots $U+2103$ relationalnormal\downarrow \downarrow $U+2193$ relationalnormal\downarrow \downarrow $U+2193$ ordinaryoperand\ee $@$ $U+2113$ ordinaryoperand\ee $@$ $U+2113$ ordinaryoperand\emptyset \emptyset $U+2205$ unaryoperand\emptyset \emptyset $U+2003$ skipnormal\end \mathbb{V} $U+3017$ closeclose\ensp $U+2025$ <th><u>\ 11 .</u></th> <th>••</th> <th>11 0000</th> <th>1.</th> <th></th>	<u>\ 11 .</u>	••	11 0000	1.	
\degree \circ U+00B0ordinaryoperand\Delta Δ U+0394ordinaryoperand\delta δ U+03B4ordinaryoperand\diamond \circ U+22C4binarynormal\diamondsuit \diamond U+2662ordinarynormal\diamondsuit \diamond U+2662ordinarynormal\diamondsuit \diamond U+2662ordinarynormal\diamondsuit \diamond U+2662ordinarynormal\dotU+00F7binarynormal\dotU+0307ordinaryaccent\dotsU+2026ordinarynormal\dotsU+2026ordinarynormal\dotsU+2103relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \downarrow U+2107ordinaryoperand\dee $@$ U+2147ordinaryoperand\ee $@$ U+213ordinaryoperand\ee $@$ U+213ordinaryoperand\end \emptyset U+2205unaryoperand\end \emptyset U+2003skipnormal\enspU+2002skipnormal\enspU+2022skipnormal\enspU+2388ordinaryoperand\enspU+2588ordinaryencl eqarray	\ddot		U+0308	ordinary	accent
(degree0+0000ordinaryoperand\Delta Δ U+0394ordinaryoperand\diamond δ U+03B4ordinaryoperand\diamondsuit δ U+22C4binarynormal\diamondsuit \diamond U+2662ordinarynormal\diamondsuit \diamond U+2662ordinarynormal\diamondsuit \diamond U+2662ordinarynormal\diamondsuit \diamond U+2662ordinarynormal\dot \because U+00F7binarynormal\dot \because U+0307ordinaryaccent\doteq \doteq U+2250relationalnormal\dotsU+2103relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \downarrow U+2193relationalnormal\downarrow \downarrow U+2193relationalnormal\downarrow \downarrow U+2193ordinaryoperand\dee e U+2147ordinaryoperand\ee e U+2147ordinaryoperand\end l U+2003skipnormal\end $]$ U+3017closeclose\enspU+2002skipnormal\ensp i U+2588ordinaryoperand\ensp i U+2588ordinaryencl eqarray			-		normal
\delta δ U+03B4ordinaryoperand\diamond \diamond U+22C4binarynormal\diamondsuit \diamond U+2662ordinarynormal\div \div U+00F7binarynormal\dot \because U+0307ordinaryaccent\doteq \doteq U+2250relationalnormal\dotsU+2026ordinarynormal\dotsU+2103relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \Downarrow U+2205ordinaryoperand\ee \mathscr{C} U+2113ordinaryoperand\ee \mathscr{C} U+2113ordinaryoperand\end ϑ U+2003skipnormal\end ϑ U+3017closeclose\enspU+2002skipnormal\enspU+2025ordinaryoperand\enspU+2028ordinaryoperand\enspU+2028ordinaryoperand\enspU+2085ordinaryoperand\enspU+2013skipnormal </td <td>\degree</td> <td>0</td> <td></td> <td>ordinary</td> <td>operand</td>	\degree	0		ordinary	operand
\diamond \diamond U+22C4binarynormal\diamondsuit \diamond U+2662ordinarynormal\div \div U+00F7binarynormal\div \div U+00F7binarynormal\dot \cdot U+0307ordinaryaccent\doteq \doteq U+2250relationalnormal\dotsU+2026ordinarynormal\dotsU+2103relationalnormal\downarrow \Downarrow U+2193relationalnormal\downarrow \downarrow U+2193ordinaryoperand\ee $@$ $U+2113$ ordinaryoperand\end \emptyset U+2205unaryoperand\end \emptyset U+2003skipnormal\end \emptyset U+3017closeclose\enspU+2002skipnormal\epsilon ϵ U+03F5ordinaryoperand\epsilon </td <td>\Delta</td> <td></td> <td>U+0394</td> <td>ordinary</td> <td>operand</td>	\Delta		U+0394	ordinary	operand
\diamondsuit \diamond U+2662ordinarynormal\div \div U+00F7binarynormal\dot \cdot U+0307ordinaryaccent\doteq \doteq U+2250relationalnormal\dotsU+2026ordinarynormal\dotsU+2103relationalnormal\downarrowUU+2193relationalnormal\downarrowUU+2193relationalnormal\dsmashUU+2193relationalnormal\dsmashUU+2207ordinaryoperand\ee \mathscr{C} U+2147ordinaryoperand\enptyset \emptyset U+2205unaryoperand\end \mathbb{J} U+2003skipnormal\enspU+2002skipnormal\enspU+2002skipnormal\enspU+2085ordinaryoperand\enspU+2002skipnormal\enspU+2002skipnormal\enspU+2002skipnormal\enspU+2588ordinaryoperand\enspU+2588ordinaryoperand	\delta	δ	U+03B4	ordinary	operand
\div \div U+00F7binarynormal\dot $U+0307$ ordinaryaccent\doteq \doteq U+2250relationalnormal\dotsU+2026ordinarynormal\dotsU+2026ordinarynormal\downarrow \Downarrow U+21D3relationalnormal\downarrow \downarrow U+2193relationalnormal\downarrow \downarrow U+2193relationalnormal\dee $@$ U+2113ordinaryoperand\end ℓ U+2205unaryoperand\end \parallel U+2003skipnormal\end \parallel U+3017closeclose\enspU+2002skipnormal\ensp ψ U+2588ordinaryoperand\ensp ψ U+2588ordinaryencl eqarray	\diamond	٥	U+22C4	binary	normal
\dot \cdot U+0307ordinaryaccent\doteq \doteq U+2250relationalnormal\dotsU+2026ordinarynormal\Downarrow \Downarrow U+21D3relationalnormal\downarrow \downarrow U+2193relationalnormal\downarrow \downarrow U+2193relationalnormal\downarrow \downarrow U+2193relationalnormal\downarrow \downarrow U+2193relationalnormal\downarrow \downarrow U+2193ordinaryencl phantom\ee $@$ U+2147ordinaryoperand\ell ℓ U+2113ordinaryoperand\enptyset \emptyset U+2205unaryoperand\enspU+2003skipnormal\enspU+2002skipnormal\enspU+2002skipnormal\enspU+2085ordinaryoperand\enspU+2085ordinaryencl eqarray	\diamondsuit	\diamond	U+2662	ordinary	normal
\dot \doteq $0+0307$ ordinaryaccent\doteq \doteq $U+2250$ relationalnormal\dots $U+2026$ ordinarynormal\Downarrow \Downarrow $U+21D3$ relationalnormal\downarrow \downarrow $U+2193$ ordinaryoperand\ee $@$ $U+2147$ ordinaryoperand\ell ℓ $U+2113$ ordinaryoperand\emptyset \emptyset $U+2205$ unaryoperand\emptyset \emptyset $U+2003$ skipnormal\end \mathbb{I} $U+3017$ closeclose\ensp $U+2002$ skipnormal\epsilon ϵ $U+03F5$ ordinaryoperand\equiv eqarray \blacksquare $U+2588$ ordinaryencl eqarray	\div	÷	U+00F7	binary	normal
$\langle abseq \rangle$ $ b 2200 \rangle$ $ b abseq \rangle$ $ b abseq \rangle$ $\langle dots$ $U + 2026$ $ordinary$ $normal$ $\langle Downarrow$ \downarrow $U + 21D3$ $relational$ $normal$ $\langle downarrow$ \downarrow $U + 2193$ $ordinary$ $operand$ $\langle ee$ e $U + 2147$ $ordinary$ $operand$ $\langle ell$ ℓ $U + 2113$ $ordinary$ $operand$ $\langle emptyset$ \emptyset $U + 2205$ $unary$ $operand$ $\langle emsp$ $U + 2003$ $skip$ $normal$ $\langle end$ $]$ $U + 3017$ $close$ $close$ $\langle ensp$ $U + 2002$ $skip$ $normal$ $\langle eqarray$ \blacksquare $U + 2588$ $ordinary$ $operand$	\dot	•	U+0307	ordinary	accent
\backslash Downarrow \Downarrow \Downarrow \Downarrow \Downarrow \Downarrow \Downarrow \square	\doteq	÷	U+2250	relational	normal
\downarrow\lambdaU+2193relationalnormal\dsmashIU+2B07ordinaryencl phantom\ee $@$ U+2147ordinaryoperand\ell ℓ U+2113ordinaryoperand\emptyset \emptyset U+2205unaryoperand\emspU+2003skipnormal\end]U+3017closeclose\enspU+2002skipnormal\epsilon ϵ U+03F5ordinaryoperand\eqarrayIU+2588ordinaryencl eqarray	\dots		U+2026	ordinary	normal
\dsmashIU+2B07ordinaryencl phantom\ee $@$ U+2147ordinaryoperand\ell ℓ U+2113ordinaryoperand\emptyset \emptyset U+2205unaryoperand\emspU+2003skipnormal\end $]$ U+3017closeclose\enspU+2002skipnormal\epsilon ϵ U+03F5ordinaryoperand\eqarray \blacksquare U+2588ordinaryencl eqarray	\Downarrow	↓	U+21D3	relational	normal
\ee \mathscr{C} U+2147ordinaryoperand\ell $\mathscr{\ell}$ U+2147ordinaryoperand\emptyset \emptyset U+2205unaryoperand\emspU+2003skipnormal\end \mathbb{J} U+3017closeclose\enspU+2002skipnormal\epsilon ϵ U+03F5ordinaryoperand\eqarray \blacksquare U+2588ordinaryencl eqarray	\downarrow	\downarrow	U+2193	relational	normal
$\[\] ell \]$ $\[\] \ell \]$ $\[\] U+2113 \]$ $\[\] ordinary \]$ $\[\] operand \]$ $\[\] emptyset \]$ $\[\] \emptyset \]$ $\[\] U+2205 \]$ $\[\] unary \]$ $\[\] operand \]$ $\[\] emsp \]$ $\[\] U+2003 \]$ $\[\] skip \]$ $\[\] normal \]$ $\[\] end \]$ $\[\] U+3017 \]$ $\[\] close \]$ $\[\] close \]$ $\[\] ensp \]$ $\[\] U+2002 \]$ $\[\] skip \]$ $\[\] normal \]$ $\[\] ensp \]$ $\[\] U+2002 \]$ $\[\] skip \]$ $\[\] normal \]$ $\[\] epsilon \]$ $\[\] \epsilon \]$ $\[\] U+03F5 \]$ $\[\] ordinary \]$ $\[\] operand \]$ $\[\] eqarray \]$ $\[\] U+2588 \]$ $\[\] ordinary \]$ $\[\] enclequarray \]$ $\[\] enclequarray \]$	\dsmash	ţ	U+2B07	ordinary	encl phantom
\backslash emptyset \emptyset $U+2205$ unaryoperand \backslash emsp $U+2003$ skipnormal \backslash end $]$ $U+3017$ closeclose \backslash ensp $U+2002$ skipnormal \backslash epsilon ϵ $U+03F5$ ordinaryoperand \backslash eqarray \blacksquare $U+2588$ ordinaryencl eqarray	\ee	e	U+2147	ordinary	operand
$\langle emsp$ $U+2003$ skipnormal $\langle end$ $]$ $U+3017$ closeclose $\langle ensp$ $U+2002$ skipnormal $\langle epsilon$ ϵ $U+03F5$ ordinaryoperand $\langle eqarray$ \blacksquare $U+2588$ ordinaryencl eqarray	\ell	ł	U+2113	ordinary	operand
\backslash end \rrbracket $U+3017$ closeclose \backslash ensp $U+2002$ skipnormal \backslash epsilon ϵ $U+03F5$ ordinaryoperand \backslash eqarray \blacksquare $U+2588$ ordinaryencl eqarray	\emptyset	Ø	U+2205	unary	operand
\enspU+2002skipnormal\epsilon ϵ U+03F5ordinaryoperand\eqarray \blacksquare U+2588ordinaryencl eqarray	\emsp		U+2003	skip	normal
\epsilon€U+03F5ordinaryoperand\eqarray■U+2588ordinaryencl eqarray	\end]	U+3017	close	close
\eqarray ■ U+2588 ordinary encl eqarray	\ensp		U+2002	skip	normal
	\epsilon	e	U+03F5	ordinary	operand
\eqno # U+0023 ordinary marker	\eqarray		U+2588	ordinary	encl eqarray
	\eqno	#	U+0023	ordinary	marker
\equiv \equiv U+2261 relational normal	\equiv	≡	U+2261	relational	normal
\eta η U+03B7 ordinary operand	\eta	η	U+03B7	ordinary	operand
\exists	\exists	Э	U+2203	unary	normal
\forall ∀ U+2200 unary normal	\forall	\forall	U+2200	unary	normal
\funcapply II U+2061 binary subsupFA	\funcapply	<i>f</i> ()	U+2061	binary	subsupFA
\Gamma Γ U+0393 ordinary operand	\Gamma	Г	U+0393	ordinary	operand
\gamma γ U+03B3 ordinary operand	\gamma	γ	U+03B3	ordinary	operand
\ge ≥ $U+2265$ relational normal	\ge	2	U+2265	relational	normal
\geq \geq U+2265 relational normal	\geq	2	U+2265	relational	normal
\gets ← U+2190 ordinary stretch horiz	\gets	\leftarrow	U+2190	ordinary	stretch horiz
\gg ≫ U+226B relational normal	\gg	>>	U+226B	relational	normal
\gimel ک U+2137 ordinary operand		ג	U+2137	ordinary	operand

\hairspU+200Askipnormal\hat^U+0302ordinaryaccent\hbarħU+210Fordinaryoperand\heartsuit♡U+2661ordinarynormal\hookleftarrow↔U+21A9relationalstretch horiz\hookrightarrow↔U+21AArelationalstretch horiz\hphantom⇔U+2B04ordinaryencl phantom	\ grave		11.0200	andinam	accort
Nat $$ U+0302ordinary ordinaryaccent\hbar \hbar U+210Fordinaryoperand\heartsuit \heartsuit U+2661ordinarynormal\hookleftarrow \leftrightarrow U+21A9relationalstretch horiz\hookrightarrow \leftrightarrow U+21A4relationalstretch horiz\hookrightarrow \leftrightarrow U+2101ordinaryencl phantom\hsmash \leftrightarrow U+2200ordinaryencl phantom\hsmash \leftrightarrow U+2201ordinaryaccent\ii l U+2148ordinaryoperand\iiint $ffff$ U+2200ordinarynary\iiint $ffff$ U+222Dordinarynary\iint $ffff$ U+2211ordinaryoperand\iint fff U+222Cordinarynary\iint fff U+2211ordinaryoperand\intn ff U+2208relationalnormal\int ff U+2208ordinaryoperand\int ff U+22149	\grave		U+0300	ordinary	accent
Nat $0+0.302$ ordinaryaccent\hbar \hbar U+210Fordinaryoperand\heartsuit \heartsuit U+2661ordinarynormal\hookleftarrow \leftrightarrow U+21A9relationalstretch horiz\hookrightarrow \leftrightarrow U+21A4relationalstretch horiz\hphantom \Leftrightarrow U+2B04ordinaryencl phantom\hsmash \leftrightarrow U+2B0Cordinaryencl phantom\hvec $^-$ U+20D1ordinaryaccent\ii i U+2148ordinaryoperand\iiint \iiint U+220Cordinarynary\iint \iiint U+222Dordinarynary\iint \iiint U+2202ordinarynary\iint \iiint U+2206unaryoperand\inth \iiint U+2208relationalnormal\inth \iiint U+2206unaryoperand\inth \iiint U+2206unaryoperand\inth f U+2218 <td></td> <td></td> <td></td> <td>-</td> <td></td>				-	
Abeartsuit \heartsuit U+2661ordinarynormalAbookleftarrow \leftrightarrow U+21A9relationalstretch horizAbookleftarrow \leftrightarrow U+21A4relationalstretch horizAbookrightarrow \leftrightarrow U+21A4relationalstretch horizAbookrightarrow \leftrightarrow U+2B04ordinaryencl phantomAbookrightarrow \leftrightarrow U+2B02ordinaryencl phantomAbookrightarrow \leftrightarrow U+2B02ordinaryencl phantomAbookrightarrow \leftrightarrow U+2148ordinaryoperandAbookrightarrow i U+2148ordinaryoperandAbookrightarrow i U+2148ordinaryoperandAbookrightarrow i U+2148ordinarynaryAbookrightarrow f U+2200ordinarynaryAbookrightarrow f U+2220ordinarynaryAbookrightarrow f U+22111ordinaryoperandAbookrightarrow f U+2206unaryoperandAbookrightarrow ϕ U+2206unaryoperandAbookrightarrow ϕ U+2206unaryoperandAbookrightarrow ϕ U+2206unaryoperandAbookrightarrow ϕ U+2206unaryoperandAbookrightarrow ϕ U+2206unaryoperandAbookrightarrow ϕ U+2206unaryoperandAbookrightarrow f U+2208ordinaryoperan	•			5	
\hookleftarrow \leftrightarrow U+21A9relationalstretch horiz\hookrightarrow \hookrightarrow U+21AArelationalstretch horiz\hphantom \Leftrightarrow U+2B04ordinaryencl phantom\hsmash \leftrightarrow U+2B0Cordinaryencl phantom\hwec $-$ U+2D01ordinaryaccent\ii l U+2148ordinaryoperand\iiint $ffffffffffffffffffffffffffffffffffff$	•	ħ		-	operand
\hookrightarrow \hookrightarrow U+21AArelationalstretch horiz\hphantom \Leftrightarrow U+2B04ordinaryencl phantom\hsmash \leftrightarrow U+2B0Cordinaryencl phantom\hvec $-$ U+20D1ordinaryaccent\ii i U+2148ordinaryoperand\iiint $fffff$ U+22D0ordinarynary\iiint $fffff$ U+22D0ordinarynary\iiint $ffffffffffffffffffffffffffffffffffff$	\heartsuit	\heartsuit	U+2661	ordinary	normal
Applantom \Leftrightarrow U+2B04ordinaryencl phantom\hsmash \leftrightarrow U+2B0Cordinaryencl phantom\hvec $-$ U+20D1ordinaryaccent\ii i U+2148ordinaryoperand\iiint $j j j$ U+22D0ordinarynary\iint $j j j$ U+22D0ordinarynary\iint $j j j$ U+222Dordinarynary\iint $j j j$ U+222Cordinaryoperand\iint $j j j$ U+2111ordinaryoperand\inath1U+0131ordinaryoperand\inath1U+2206unaryoperand\infty ∞ U+221Eordinaryoperand\infty ∞ U+221Eordinaryoperand\infty ∞ U+221Eordinaryoperand\infty ∞ U+221Eordinaryoperand\infty $j j$ U+222Bordinaryoperand\infty $j j$ U+221Eordinaryoperand\int $j j$ U+221Bordinaryoperand\inta $j j$ U+2149ordinaryoperand\inta $j j$ U+2149ordinaryoperand <t< td=""><td>\hookleftarrow</td><td>ب ب</td><td>U+21A9</td><td>relational</td><td>stretch horiz</td></t<>	\hookleftarrow	ب ب	U+21A9	relational	stretch horiz
\hsmash \leftrightarrow U+2B0Cordinaryencl phantom\hvec \cdot U+20D1ordinaryaccent\ii i U+2148ordinaryoperand\iiint $ffff$ U+2240Cordinarynary\iiint $ffff$ U+222Dordinarynary\iint $ffff$ U+222Dordinarynary\iint $ffff$ U+222Cordinarynary\iint $ffffffffffffffffffffffffffffffffffff$	\hookrightarrow	\hookrightarrow	U+21AA	relational	stretch horiz
\hvecIU+20D1ordinaryaccent\ii i U+20D1ordinaryoperand\iiint j U+2148ordinarynary\iiint j U+220Cordinarynary\iint j U+222Dordinarynary\iint j U+222Cordinarynary\iint j U+2111ordinaryoperand\imath1U+0131ordinaryoperand\inath ϵ U+2208relationalnormal\inc Δ U+2206unaryoperand\infty ∞ U+221Eordinaryoperand\infty ∞ U+221Eordinaryoperand\int f U+2228ordinaryoperand\int f U+2218ordinaryoperand\int f U+2218ordinaryoperand\int f U+2219ordinaryoperand\inta f U+2149ordinaryoperand\inta f U+228ordinaryoperand\inta f U+237ordinaryoperand\inta f U+238ordinaryoperand\inta f U+039Bordinaryoperand\inta f U+039Bordinaryoperand\inta f U+27E9closeclose\Lambda A U+038Bordinaryoperand\lambda A U+038Bordinaryope	\hphantom	⇔	U+2B04	ordinary	encl phantom
\ii i U+2148ordinaryoperand\iiint \iiint U+2148ordinarynary\iiint \iiint U+2200ordinarynary\iint \iiint U+2220ordinarynary\iint \iint U+2220ordinarynary\imath \iint U+2220ordinaryoperand\imath1U+0131ordinaryoperand\inath \in U+2208relationalnormal\inc Δ U+2206unaryoperand\infty ∞ U+221Eordinaryoperand\int \int U+222Bordinarynary\iota ι U+03B9ordinaryoperand\int \int U+2149ordinaryoperand\iptica ι U+03B9ordinaryoperand\iptica j U+2149ordinaryoperand\iptica χ U+03B9ordinaryoperand\iptica χ U+03BAordinaryoperand\iptica Λ U+03BBordinaryoperand\iptica Λ U+03BBordinaryoperand\iptica χ U+27E8openopen\lambda Λ U+03BBordinaryoperand\iptica $\{$ U+007Bopenopen\lambda Λ U+03B8ordinaryoperand\lambda $\{$ U+07B8openopen\lambda $\{$ U+07Bopen <td< td=""><td>\hsmash</td><td>‡</td><td>U+2B0C</td><td>ordinary</td><td>encl phantom</td></td<>	\hsmash	‡	U+2B0C	ordinary	encl phantom
Niiint \iiint U+2A0CordinarynaryNiiint \iiint U+2A0CordinarynaryNiint \iiint U+222DordinarynaryNimt \iint U+222CordinarynaryNamath1U+222CordinaryoperandNimath1U+2111ordinaryoperandNimath1U+2108relationalnormalNinc Δ U+2206unaryoperandNint \int U+221EordinaryoperandNint \int U+222BordinaryoperandNint \int U+2149ordinaryoperandNint \int U+2149ordinaryoperandNint \int U+2149ordinaryoperandNint \int U+2149ordinaryoperandNint \int U+218ordinaryoperandNinth \int U+237ordinaryoperandNinth \int U+238ordinaryoperandNinth \int U+27E9closecloseNambha Λ U+03BBordinaryoperandNambha Λ U+03BBordinaryoperandNambha Λ U+03BBordinaryoperandNambha Λ U+03BBordinaryoperandNambha Λ U+03BBopenopenNambha Λ U+03BBopenopenNambha Λ U+07BopenopenN	\hvec	-	U+20D1	ordinary	accent
\iiint \iiint U+222Dordinarynary\iint \iint U+222Cordinarynary\Im \Im U+2111ordinaryoperand\imath1U+0131ordinaryoperand\imath1U+0131ordinaryoperand\in \in U+2208relationalnormal\inc Δ U+2206unaryoperand\int \int U+221Eordinaryoperand\int \int U+221Eordinaryoperand\int \int U+222Bordinaryoperand\int \int U+222Bordinaryoperand\int \int U+221Fordinaryoperand\int \int U+221Bordinaryoperand\int \int U+221Fordinaryoperand\int \int U+221Bordinaryoperand\int \int U+221Fordinaryoperand\int \int U+2219closeclose\int J U+039Bordinaryoperand\ket \rangle U+27E9closeclose\Lambda Λ U+039Bordinaryoperand\langle \langle U+27E8openopen\langle \langle U+27E8openopen\langle $[$ U+007Bopenopen\langle $[$ U+2308openopen\langle $[$ U+2308openopen\langle<	\ii	11	U+2148	ordinary	operand
\iint \iint U+222Cordinarynary\Im \Im U+2111ordinaryoperand\imath1U+0131ordinaryoperand\in \in U+2208relationalnormal\inc Δ U+2206unaryoperand\infty ∞ U+221Eordinaryoperand\int \int U+222Bordinarynary\iota ι U+03B9ordinaryoperand\int \int U+2149ordinaryoperand\jj j U+2149ordinaryoperand\jmathJU+03B9ordinaryoperand\kappa κ U+03BAordinaryoperand\ket \rangle U+27E9closeclose\Lambda Λ U+03BBordinaryoperand\langle(U+27E8openopen\langle \langle U+27E8openopen\langle \langle U+2788openopen\langle \langle U+2788openopen\langle \langle U+207Bopenopen\langle \langle U+2308openopen\langle \int U+2308openopen\langle \int U+2215binarydivide\langle \int U+2215ordinarynormal	\iiiint		U+2A0C	ordinary	nary
\Im \Im U+2111ordinaryoperand\imath1U+0131ordinaryoperand\in \in U+2208relationalnormal\inc Δ U+2206unaryoperand\infty ∞ U+221Eordinaryoperand\int \int U+222Bordinarynary\iota1U+03B9ordinaryoperand\ijj j U+2149ordinaryoperand\jipathJU+03B9ordinaryoperand\jmathJU+0237ordinaryoperand\kappa κ U+03BAordinaryoperand\kappa Λ U+03BBordinaryoperand\lambda Λ U+03BBopenopen\lambda Λ U+07Bopenopen\lambda $[$ U+07Bopenopen\lambda $[$ U+007Bopenopen\lambda $[$ U+2308openopen\lambda $[$ U+2308openopen\lambda $[$ U+2215binarydivide	\iiint		U+222D	ordinary	nary
\imath1U+0131ordinaryoperand\in \in U+2208relationalnormal\inc Δ U+2206unaryoperand\infty ∞ U+221Eordinaryoperand\int \int U+222Bordinarynary\iota ι U+03B9ordinaryoperand\jji j U+2149ordinaryoperand\jipathJU+0237ordinaryoperand\kappa κ U+03BAordinaryoperand\kappa κ U+03BAordinaryoperand\kath λ U+03BAordinaryoperand\lambda Λ U+03BBordinaryoperand\lambda Λ U+03BBopenopen\lambda \langle U+27E8openopen\lambda $[$ U+007Bopenopen\lambda $[$ U+2308<	\iint	∬	U+222C	ordinary	nary
$\langle in $ \in U+2208relationalnormal $\langle inc $ Δ U+2206unaryoperand $\langle infty $ ∞ U+221Eordinaryoperand $\langle int $ \int U+222Bordinarynary $\langle iota $ ι U+03B9ordinaryoperand $\langle ijj $ j' U+2149ordinaryoperand $\langle ipi $ ipi U+2149ordinaryoperand $\langle ipi $ j' U+2149ordinaryoperand $\langle ipi $ j' U+27E9closeclose $\langle kappa $ κ U+03BAordinaryoperand $\langle ket $ \rangle U+27E9closeclose $\langle Lambda $ Λ U+03BBordinaryoperand $\langle langle $ \langle U+27E8openopen $\langle lbrace $ $\{$ U+007Bopenopen $\langle lordk $ $[$ U+005Bopenopen $\langle ldiv $ $/$ U+2308openopen $\langle ldis $ \ldots U+2215binarydivide $\langle ldots$ \ldots U+2026ordinarynormal	\Im	II	U+2111	ordinary	operand
$eq:linear_li$	\imath	1	U+0131	ordinary	operand
\infty ∞ U+221Eordinaryoperand\int \int U+222Bordinarynary\iota ι U+03B9ordinaryoperand\jj j U+2149ordinaryoperand\jmathJU+2149ordinaryoperand\kappa κ U+03BAordinaryoperand\ket \rangle U+27E9closeclose\Lambda Λ U+039Bordinaryoperand\langle \langle U+27E9closeclose\langle \langle U+039Bordinaryoperand\langle \langle U+27E8openopen\lbrace $\{$ U+007Bopenopen\lbrack[U+007Bopenopen\ldiv/U+2215binarydivide\ldotsU+2206ordinarynormal	\in	E	U+2208	relational	normal
\int \int U+222Bordinarynary\iota ι U+03B9ordinaryoperand\jj j U+2149ordinaryoperand\jmath j U+0237ordinaryoperand\kappa κ U+03BAordinaryoperand\ket \rangle U+27E9closeclose\Lambda Λ U+03BBordinaryoperand\lambda λ U+03BBordinaryoperand\lambda \langle U+27E9closeclose\lambda \langle U+03BBordinaryoperand\lambda \langle U+03BBordinaryoperand\lambda \langle U+03BBordinaryoperand\lambda \langle U+03BBordinaryoperand\lambda \langle U+27E8openopen\lambda \langle U+27E8openopen\langle \langle U+205Bopenopen\latic $[$ U+005Bopenopen\latic $[$ U+2308openopen\latic $[$ U+2215binarydivide\laticU+2226ordinarynormal	\inc	Δ	U+2206	unary	operand
\iotaιU+03B9ordinaryoperand\jjjU+2149ordinaryoperand\jmathjU+0237ordinaryoperand\kappaκU+03BAordinaryoperand\ket>U+27E9closeclose\LambdaΛU+039Bordinaryoperand\lambdaΛU+039Bordinaryoperand\lambda(U+039Bordinaryoperand\lambda(U+039Bordinaryoperand\lambda(U+039Bordinaryoperand\lambda(U+039Bordinaryoperand\lambda(U+039Bordinaryoperand\lambda(U+039Bopenopen\lambda(U+27E8openopen\lambda(U+27E8openopen\langle(U+2308openopen\langle[U+2308openopen\langle(U+2215binarydivide\langleU+2215binarynormal	\infty	∞	U+221E	ordinary	operand
\jjjU+2149ordinaryoperand\jmathJU+0237ordinaryoperand\kappaκU+03BAordinaryoperand\ket>U+27E9closeclose\LambdaΛU+039Bordinaryoperand\lambdaΛU+039Bordinaryoperand\lambda(U+039Bordinaryoperand\lambda(U+039Bordinaryoperand\lambda(U+03BBordinaryoperand\lambda(U+03BBordinaryoperand\langle(U+03BBopenopenU+007Bopenopen\lbrack[U+005Bopenopen\lceil[U+2308openopen\ldiv/U+2215binarydivide\ldotsU+2026ordinarynormal	\int	ſ	U+222B	ordinary	nary
\jmathJU+0237ordinaryoperand\kappaκU+03BAordinaryoperand\ket>U+27E9closeclose\LambdaΛU+039Bordinaryoperand\lambdaλU+03BBordinaryoperand\langle<	\iota	ι	U+03B9	ordinary	operand
\kappaκU+03BAordinaryoperand\ket>U+27E9closeclose\LambdaΛU+039Bordinaryoperand\lambdaλU+03BBordinaryoperand\langle<	\jj	Ĵ	U+2149	ordinary	operand
\ket>U+27E9closeclose\LambdaΛU+039Bordinaryoperand\lambdaλU+03BBordinaryoperand\langle<	\jmath	J	U+0237	ordinary	operand
\LambdaΛU+039Bordinaryoperand\lambdaλU+03BBordinaryoperand\langle<	\kappa	κ	U+03BA	ordinary	operand
\lambdaλU+03BBordinaryoperand\langle<	\ket	>	U+27E9	close	close
\langle(U+27E8openopenU+007Bopenopen\lbrack[U+005Bopenopen\lceil[U+2308openopen\ldiv/U+2215binarydivide\ldotsU+2026ordinarynormal	\Lambda	Λ	U+039B	ordinary	operand
U+007Bopenopen\lbrack[U+005Bopenopen\lceil[U+2308openopen\ldiv/U+2215binarydivide\ldotsU+2026ordinarynormal	\lambda	λ	U+03BB	ordinary	operand
\lbrack[U+005Bopenopen\lceil[U+2308openopen\ldiv/U+2215binarydivide\ldotsU+2026ordinarynormal	\langle	<	U+27E8	open	open
\lbrack[U+005Bopenopen\lceil[U+2308openopen\ldiv/U+2215binarydivide\ldotsU+2026ordinarynormal	\lbrace	{	U+007B	open	open
\lceil[U+2308openopen\ldiv/U+2215binarydivide\ldotsU+2026ordinarynormal	\lbrack	[U+005B	open	open
\ldots U+2026 ordinary normal	\lceil	ſ	U+2308	open	open
	\ldiv	/	U+2215	binary	divide
$le \leq U+2264$ relational normal	\ldots		U+2026	ordinary	normal
	\le	≤	U+2264	relational	normal
\Leftarrow \leftarrow U+21D0 relational stretch horiz	\Leftarrow	⇐	U+21D0	relational	stretch horiz

Unicode Technical Note 28

Vertication $-$ U+21BDrelationalstretch horizVeftharpoonup $-$ U+21BCrelationalstretch horizVeftrightarrow \leftrightarrow U+21D4relationalstretch horizVeftrightarrow \leftrightarrow U+2194relationalstretch horizVeftrightarrow \leftrightarrow U+2264relationalnormalVeftrightarrow \leftarrow U+2264relationalnormalVeftrightarrow \leftarrow U+27F8relationalnormalVongleftarrow \leftarrow U+27F7relationalnormalVongleftrightarrow \leftarrow U+27F7relationalnormalVongleftrightarrow \leftrightarrow U+27F7relationalnormalVongleftrightarrow \leftrightarrow U+27F7relationalnormalVongrightarrow \rightarrow U+27F6relationalnormalNapsto \rightarrow U+27F6relationalnormalNampsto \rightarrow U+225A0ordinaryencl matrixNedspU+225FOrdinarynormalNindels \models U+2213unary/binaryNundels \models U+2207unaryNapsto \checkmark U+2292ordinaryoperandNapsp $-$ U+2298binarydivideNapsp $-$ U+2298binarydivideNapsp $-$ U+2206relationalnormalNapsp $-$ U+2206relationalnormalNapsp $-$ U+2208relationalnormal <t< th=""><th></th><th></th><th></th><th></th><th>1</th></t<>					1
Neffharpoonup \leftarrow U+21BCrelationalstretch horizLeftrightarrow \leftrightarrow U+21D4relationalstretch horizNeffrightarrow \leftrightarrow U+2194relationalstretch horizNeffrightarrow \leftrightarrow U+2194relationalstretch horizNeffrightarrow \leftarrow U+2264relationalnormalNeloco $[$ U+226ArelationalnormalNeloco $($ U+226ArelationalnormalNongleftarrow \leftarrow U+27F8relationalnormalNongleftarrow \leftarrow U+27F7relationalnormalNongleftrightarrow \leftrightarrow U+27F7relationalnormalNongrightarrow \leftrightarrow U+27F6relationalnormalNongrightarrow \rightarrow U+27F6relationalnormalNapsto \mapsto U+27F6relationalstretch horizNatrix \blacksquare U+225Aordinaryencl matrixNedsp \cup U+2275OrdinarynormalNind1U+2233relationalstretch horizNmid1U+2234relationalstretch horizNmp \mp U+2205OrdinaryoperandNabla ∇ U+2207unary/binaryunary/binaryNull μ U+038CordinaryoperandNabla ∇ U+2298binarydivideNearrow \wedge U+2207relationalnormalNearrow \wedge U+2208 <t< td=""><td>\leftarrow</td><td>\leftarrow</td><td>U+2190</td><td>relational</td><td>stretch horiz</td></t<>	\leftarrow	\leftarrow	U+2190	relational	stretch horiz
Leftrightarrow \Leftrightarrow U+21D4relationalstretch horiz\leftrightarrow \leftrightarrow U+2194relationalstretch horiz\leq \leq U+2264relationalnormal\lfloor $ $ U+230Aopenopen\ll \ll U+226Arelationalnormal\longleftarrow \leftarrow U+27F8relationalnormal\longleftrightarrow \leftarrow U+27F7relationalnormal\longleftrightarrow \leftarrow U+27F7relationalnormal\longrightarrow \leftrightarrow U+27F7relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalstretch horiz\mathbf{matrix \blacksquare U+2250ordinaryencl matrix\mathbf{matrix \blacksquare U+2265Ordinarynormal\mathbf{matrix \blacksquare U+2207unary/binaryunary/binary\mathbf{matrix \blacksquare U+2207unary/binaryoperand\mathbf{matrix \bigcirc U+2292ordinaryoperand\mathbf{matrix} \bigcirc U+2292ordinaryoperand\mathbf{matrix} \bigcirc U+2207unaryoperand\mathbf{matrix} \bigcirc U+2292ordinaryoperand\mathbf{matrix} \bigcirc U+2292ordinaryoperand\mathbf{matrix} \square U+2293inaryidi	\leftharpoondown	Ţ	U+21BD	relational	stretch horiz
\leftrightarrow \leftrightarrow U+2194relationalstretch horiz\leq \leq U+2264relationalnormal\lfloor U+230Aopenopen\llU+230Arelationalnormal\longleftarrow \leftarrow U+27F8relationalnormal\longleftarrow \leftarrow U+27F5relationalnormal\longleftrightarrow \leftarrow U+27F7relationalnormal\longleftrightarrow \leftarrow U+27F7relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalstretch horiz\mathbf{matrix \blacksquare U+225Aordinaryencl matrix\mathbf{matrix \blacksquare U+2273relationalstretch horz\mathbf{matrix \blacksquare U+2273relationalstretch horz\mathbf{matrix \blacksquare U+2233relationalstretch horz\mathbf{matrix \blacksquare U+2248relationalstretch horz\mathbf{matrix} \blacksquare U+2273unary/binaryunary/binary\mathbf{matrix} \blacksquare U+2233relationalstretch horz\mathbf{matrix} \blacksquare U+2248relationalstretch horz\mathbf{matrix} \blacksquare U+2248relationalstretch horz\mathbf{matrix} \blacksquare U+22	\leftharpoonup	<u> </u>	U+21BC	relational	stretch horiz
\leq \leq U+2264relationalnormal\lfloor U+230Aopenopen\ll \ll U+230Aopenopen\longleftarrow \leftarrow U+27F8relationalnormal\longleftarrow \leftarrow U+27F5relationalnormal\longleftrightarrow \leftarrow U+27F7relationalnormal\longleftrightarrow \leftrightarrow U+27F7relationalnormal\longleftrightarrow \rightarrow U+27F7relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalstretch horiz\mathbf{matrix \blacksquare U+25A0ordinaryencl matrix\mathbf{matrix \blacksquare U+2258Ordinarynormal\mathbf{matrix \blacksquare U+2233relationalstretch horiz\mathbf{matrix \blacksquare U+2248relationalstretch horiz\mathbf{matrix \blacksquare U+2233relationalstretch horiz\mathbf{matrix \blacksquare U+2248relationalstretch horiz\mathbf{matrix} \blacksquare <t< td=""><td>\Leftrightarrow</td><td>\Leftrightarrow</td><td>U+21D4</td><td>relational</td><td>stretch horiz</td></t<>	\Leftrightarrow	\Leftrightarrow	U+21D4	relational	stretch horiz
NinoIU+230Aopenopen\ll $<$ U+230Arelationalnormal\longleftarrow \leftarrow U+27F8relationalnormal\longleftarrow \leftarrow U+27F5relationalnormal\longleftrightarrow \leftarrow U+27F7relationalnormal\longleftrightarrow \leftrightarrow U+27F7relationalnormal\longrightarrow \leftrightarrow U+27F7relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\mapsto \mapsto U+27F6relationalstretch horiz\mathbf{matrix \blacksquare U+223relationalstretch horiz\mathbf{matrix \blacksquare U+2248relationalstretch horz\mathbf{matrix \blacksquare U+2248relationalstretch horz\mathbf{matrix \blacksquare U+2213unary/binaryunary/binary\mathbf{matrix \blacksquare U+2207unaryoperand\mathbf{matrix \square U+2392ordinaryoperand\mathbf{matrix} \square U+2393inaryidvide\mathbf{matrix} \square U+2207unaryoperand\mathbf{matrix} \square U+2393idiandnormal\mathbf{matrix} \square U+2394inaryidivide\mathbf{matrix} \square U+2394relationalnormal <td>\leftrightarrow</td> <td>\leftrightarrow</td> <td>U+2194</td> <td>relational</td> <td>stretch horiz</td>	\leftrightarrow	\leftrightarrow	U+2194	relational	stretch horiz
Image: Normal intermediateImage: Normal intermediate\ll $<$ U+226Arelationalnormal\Longleftarrow \leftarrow U+27F8relationalnormal\Longleftrightarrow \leftarrow U+27F7relationalnormal\Longrightarrow \leftrightarrow U+27F7relationalnormal\Longrightarrow \leftrightarrow U+27F6relationalnormal\Longrightarrow \rightarrow U+27F6relationalnormal\Longrightarrow \rightarrow U+27F6relationalnormal\mapsto \rightarrow U+27F6relationalstretch horiz\mapsto \rightarrow U+27F6relationalstretch horiz\mapsto \rightarrow U+25A0ordinaryencl matrix\mapsto \cup U+205FOrdinarynormal\mapsto \vdash U+2233relationalstretch horz\mapsto \vdash U+2248relationalstretch horz\mapsto \vdash U+2213unary/binaryunary/binary\mapsto \vdash U+2213unary/binaryunary/binary\mapsto \downarrow U+2207unaryoperand\mapsto \downarrow U+2207unaryoperand\mapsto \downarrow U+2208relationalnormal\mapsto \downarrow U+2209inaryoperand\mapsto \downarrow U+2209relationalnormal\mapsto \downarrow U+2208relationalnormal\mapsto \downarrow U+2208relationalnorm	\leq	\leq	U+2264	relational	normal
\Longleftarrow \leftarrow U+27F8relationalnormal\longleftarrow \leftarrow U+27F5relationalnormal\Longleftrightarrow \leftrightarrow U+27F7relationalnormal\longrightarrow \rightarrow U+27F7relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\mapsto \rightarrow U+27F6relationalstretch horiz\mapsto \rightarrow U+21A6relationalstretch horiz\mathbf{matrix} \blacksquare U+205FOrdinarynormal\mathbf{mdsp} \square U+205FOrdinarynormal\mathbf{mdsp} \square U+223relationallist delims\mathbf{mdsp} \square U+2213unary/binaryunary/binary\mathbf{mdls} \models U+2207unaryoperand\mathbf{mdls} \neg U+2592ordinaryoperand\nabla ∇ U+2207unaryoperand\nabla ∇ U+2207unaryoperand\nabla ∇ U+2208binarydivide\ne \neq U+2260relationalnormal\nhearrow \land U+2197relationalnormal\ne \neq U+2208binarydivide\ne \neq U+2208relationalnormal\neq \neq U+2208relationalnormal\neq \downarrow U+2208relationalnorma	\lfloor	l	U+230A	open	open
\longleftarrow←U+27F5relationalnormal\Longleftrightarrow⇔U+27F7relationalnormal\longleftrightarrow⇒U+27F7relationalnormal\Longrightarrow⇒U+27F6relationalnormal\longrightarrow→U+27F6relationalnormal\mapsto↦U+27F6relationalstretch horizmatrix■U+25A0ordinaryencl matrix\medspU+205FOrdinarynormal\midIU+2223relationallist delims\models⊨U+2213unary/binaryunary/binary\muµU+03BCordinaryoperand\nabla∇U+2207unaryoperand\nabla∇U+2298binarydivide\ne≠U+2208relationalnormal\nhiv∅U+2298binarydivide\ne≠U+2207relationalnormal\nearrow∧U+2197relationalnormal\nearrow√U+2298binarydivide\ne≠U+2260relationalnormal\neq↓U+2260relationalnormal\neq↓U+2260relationalnormal\neq↓U+2260relationalnormal\neq↓U+2260relationalnormal\neq↓U+2260relationalnormal	\ll	~	U+226A	relational	normal
\Longleftrightarrow⇔U+27FArelationalnormal\longleftrightarrow↔U+27F7relationalnormal\Longrightarrow→U+27F6relationalnormal\longrightarrow→U+27F6relationalnormal\mapsto↦U+27F6relationalstretch horiz\mapsto↦U+25A0ordinaryencl matrix\medspU+205FOrdinarynormal\middels⊨U+2223relationallist delims\models⊨U+2213unary/binaryunary/binary\muµU+03BCordinaryoperand\nabla∇U+2207unaryoperand\nablaVU+2298binarydivide\nearrow/2U+2197relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208relationalnormal\nearrow/2U+2208 <td>\Longleftarrow</td> <td>\Downarrow</td> <td>U+27F8</td> <td>relational</td> <td>normal</td>	\Longleftarrow	\Downarrow	U+27F8	relational	normal
Normal\longleftrightarrow \leftrightarrow U+27F7relationalnormal\Longrightarrow \rightarrow U+27F6relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\mapsto \mapsto U+21A6relationalstretch horiz\matrix \blacksquare U+25A0ordinaryencl matrix\medsp $U+205F$ Ordinarynormal\middIU+2223relationallist delims\models \models U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\nabla ∇ U+2292ordinarynormal\nbsp $U+00A0$ skipnormal\nliv \mathcal{O} U+2197relationalnormal\nearrow 2 U+2197relationalnormal\neq \neq U+2208relationalnormal\neq ϕ U+2208relati	\longleftarrow	\leftarrow	U+27F5	relational	normal
\Longrightarrow \Rightarrow U+27F9relationalnormal\longrightarrow \rightarrow U+27F6relationalnormal\mapsto \mapsto U+21A6relationalstretch horiz\matrix \blacksquare U+25A0ordinaryencl matrix\medspU+205FOrdinarynormal\midIU+223relationallist delims\models \models U+2248relationalstretch horz\mp \mp U+2213unary/binaryunary/binary\mu<	\Longleftrightarrow	€	U+27FA	relational	normal
\longrightarrow \rightarrow U+27F6relationalnormal\mapsto \mapsto U+21A6relationalstretch horiz\matrix \blacksquare U+25A0ordinaryencl matrix\medspU+205FOrdinarynormal\midIU+2223relationallist delims\models \models U+22A8relationalstretch horz\mp \mp U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\nabla ∇ U+2292ordinarynormal\nhispU+00A0skipnormal\ndiv \oslash U+2298binarydivide\nearrow \checkmark U+2197relationalnormal\neq \neq U+2206relationalnormal\neq \Rightarrow U+2260relationalnormal\neq \neg U+2207unarynormal\neq \checkmark U+2298binarydivide\neq \neg U+2260relationalnormal\neq \neg U+2206relationalnormal\neq \neg U+2208relationalnormal\neq \neg U+2208relationalnormal\neq \neg U+2208relationalnormal\neq \neg U+2208relationalnormal\ni \ni U+2208relationalnormal\norm $ $ U+21	\longleftrightarrow	\leftrightarrow	U+27F7	relational	normal
\mapsto \mapsto U+21A6relationalstretch horiz\matrix \blacksquare U+21A6ordinaryencl matrix\medspU+205FOrdinarynormal\midIU+205FOrdinarynormal\midIU+223relationallist delims\models \models U+22A8relationalstretch horz\mp \mp U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\nabla ∇ U+2292ordinaryoperand\nbspU+00A0skipnormal\ndiv \mathcal{O} U+2298binarydivide\nearrow \mathcal{A} U+2197relationalnormal\neq \neq U+2206relationalnormal\neq \neq U+2208relationalnormal\neq \neq U+2208relationalnormal\neq \neq U+2208relationalnormal\norm \parallel U+2016ordinaryopen/close\nu ν U+2196relationalnormal\norm \parallel U+2196relationalnormal\norm \square U+2196relationalnormal\norm \square U+2208relationalnormal\norm \square U+2196relationalnormal\norm \square U+2196relationalnormal\norm \square <td< td=""><td>\Longrightarrow</td><td>\Rightarrow</td><td>U+27F9</td><td>relational</td><td>normal</td></td<>	\Longrightarrow	\Rightarrow	U+27F9	relational	normal
NatrixU+25A0ordinaryencl matrix\madspU+205FOrdinarynormal\midIU+205FOrdinarynormal\midIU+2223relationallist delims\models \models U+22A8relationalstretch horz\mp \mp U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\naryandIU+2592ordinarynormal\nbspU+0A0skipnormal\ndiv \oslash U+2208binarydivide\nearrow \checkmark U+2197relationalnormal\neg \neg U+00ACunarynormal\neg \neg U+2208relationalnormal\neg \neg U+2208relationalnormal\norm \parallel U+2208relationalnormal\norm \Downarrow U+2208relationalnormal\norm \checkmark U+2208relationalnormal\norm \Downarrow U+2016ordinaryopen/close\nu ν U+33BDordinaryoperand\nwarrow \searrow U+2196relationalnormal\nodot \bigcirc U+2208inaryoperand	\longrightarrow	\rightarrow	U+27F6	relational	normal
\medspU+205FOrdinarynormal\midIU+205FOrdinarynormal\midIU+2223relationallist delims\models \models U+22A8relationalstretch horz\mp \mp U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\nabla ∇ U+2207unaryoperand\naryand \blacksquare U+2592ordinarynormal\nbspU+00A0skipnormal\ndiv \oslash U+2298binarydivide\ne \neq U+2260relationalnormal\nearrow \land U+2197relationalnormal\neq \neq U+2208relationalnormal\neq \neq U+2208relationalnormal\ni \ni U+2208relationalnormal\norm \parallel U+2016ordinaryoperand\nu ν U+03BDordinaryoperand\nu ψ U+2196relationalnormal	\mapsto	→	U+21A6	relational	stretch horiz
\midIU+2223relationallist delims\models \models U+22A8relationalstretch horz\mp \mp U+22A8relationalstretch horz\mp \mp U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\naryand $rac{1}{2}$ U+2592ordinarynormal\nbspU+00A0skipnormal\ndiv O U+2298binarydivide\ne \neq U+2260relationalnormal\nearrow \checkmark U+2197relationalnormal\neq \neq U+2208relationalnormal\neq \neq U+2208relationalnormal\ni \ni U+2016ordinaryopen/close\nu ν U+03BDordinaryoperand\nwarrow \land U+2196relationalnormal\nwarrow \heartsuit U+2208relationalnormal\norm \parallel U+2016ordinaryopen/close\nu ν U+03BDordinaryoperand	\matrix		U+25A0	ordinary	encl matrix
\models \models U+22A8relationalstretch horz\mp \mp U+22A8relationalstretch horz\mp \mp U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\naryand im U+2592ordinarynormal\nbspU+00A0skipnormal\ndiv \oslash U+2298binarydivide\ne \neq U+2260relationalnormal\nearrow \checkmark U+2197relationalnormal\neg \neg U+00ACunarynormal\neg \neg U+2260relationalnormal\neg \neg U+2260relationalnormal\neg \neg U+2208relationalnormal\ni \ni U+2208relationalnormal\norm \parallel U+2016ordinaryopen/close\nu ν U+03BDordinaryoperand\norm \parallel U+2196relationalnormal\norm \neg U+2298binaryoperand	\medsp		U+205F	Ordinary	normal
\mp $\overline{\mp}$ U+2213unary/binaryunary/binary\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\naryand \overline{w} U+2592ordinarynormal\nbspU+00A0skipnormal\ndiv \oslash U+2298binarydivide\ne \neq U+2260relationalnormal\nearrow \checkmark U+2197relationalnormal\neg \neg U+00ACunarynormal\neg \neg U+2260relationalnormal\neg \neg U+2208relationalnormal\norm \parallel U+2208relationalnormal\norm \downarrow U+2166ordinaryopen/close\nu ν U+03BDordinaryoperand\nwarrow \nwarrow U+2196relationalnormal\normal \bigcirc U+2299binarynormal	\mid	I	U+2223	relational	list delims
\mu μ U+03BCordinaryoperand\nabla ∇ U+2207unaryoperand\naryand \blacksquare U+2592ordinarynormal\nbspU+00A0skipnormal\nbspU+0298binarydivide\ne \neq U+2260relationalnormal\nearrow \wedge U+2197relationalnormal\neg \neg U+00ACunarynormal\neg \neg U+2208relationalnormal\neq \neq U+2208relationalnormal\ni \ni U+2208relationalnormal\norm \parallel U+2016ordinaryopen/close\nu ν U+03BDordinaryoperand\nwarrow \ddots U+2196relationalnormal\odot \odot U+2299binarynormal	\models	Ш	U+22A8	relational	stretch horz
∇ ∇ U+2207unaryoperand \naryand \blacksquare U+2592ordinarynormal \nbsp U+00A0skipnormal \ndiv \oslash U+2298binarydivide \ne \neq U+2260relationalnormal \ne \neq U+2197relationalnormal \neg \neg U+00ACunarynormal \neg \neg U+2260relationalnormal \neg \neg U+2208relationalnormal \neq \neq U+2208relationalnormal \norm \parallel U+2016ordinaryopen/close ν ν U+03BDordinaryoperand \nwarrow \normal \normal normal \normal \normal \normal normal ν ν U+2196relationalnormal \normal \normal \normal normal \normal \normal \normal normal ν	\mp	Ŧ	U+2213	unary/binary	unary/binary
\naryand \blacksquare U+2592ordinarynormal\nbspU+00A0skipnormal\ndiv \oslash U+2298binarydivide\ne \neq U+2260relationalnormal\ne \neq U+2197relationalnormal\neg \neg U+00ACunarynormal\neg \neg U+00ACunarynormal\neg \neg U+2200relationalnormal\neq \neq U+2208relationalnormal\ni \ni U+220Brelationalnormal\norm \parallel U+2016ordinaryopen/close\nu ν U+03BDordinaryoperand\nwarrow \searrow U+2196relationalnormal\odot \odot U+2299binarynormal	\mu	μ	U+03BC	ordinary	operand
\nbspU+00A0skipnormal\ndiv∅U+2298binarydivide\ne≠U+2260relationalnormal\nearrow↗U+2197relationalnormal\neg¬U+00ACunarynormal\neq≠U+2260relationalnormal\neq≠U+2260relationalnormal\niЭU+220Brelationalnormal\niYU+2016ordinaryopen/close\nu∨U+03BDordinaryoperand\nwarrow≦U+2196relationalnormal\odot⊙U+2299binarynormal	\nabla	∇	U+2207	unary	operand
\ndiv \oslash U+2298binarydivide\ne \neq U+2260relationalnormal\nearrow \checkmark U+2197relationalnormal\neg \neg U+00ACunarynormal\neg \neg U+2260relationalnormal\neq \neq U+220Brelationalnormal\ni \ni U+220Brelationalnormal\norm \parallel U+2016ordinaryopen/close\nu ν U+2196relationalnormal\nwarrow \nwarrow U+2196relationalnormal\odot \odot U+2299binarynormal	\naryand		U+2592	ordinary	normal
$\[\] ne$ \neq U+2260relationalnormal $\[\] nearrow$ $\[\] normal$ $\[\] u+2197$ relationalnormal $\[\] neg$ $\[\] normal$ $\[\] u+200AC$ unarynormal $\[\] neq$ $\[\] u+2260$ relationalnormal $\[\] neq$ $\[\] u+220B$ relationalnormal $\[\] ni$ $\[\] u+220B$ relationalnormal $\[\] norm$ $\[\] u+2016$ ordinaryopen/close $\[\] nu$ $\[\] v$ $\[u+2196$ relationalnormal $\[\] nwarrow$ $\[\] v$ $\[u+2196$ relationalnormal $\[\] odot$ $\[\] output\[u+2299binarynormal$	\nbsp		U+00A0	skip	normal
\nearrowU+2197relationalnormal\neg¬U+00ACunarynormal\neq≠U+2260relationalnormal\ni∋U+220Brelationalnormal\norm U+2016ordinaryopen/close\nu∨U+03BDordinaryoperand\nwarrow\U+2196relationalnormal\odot⊙U+2299binarynormal	\ndiv	\oslash	U+2298	binary	divide
$\[\] \[\] \[\] \[\] \[\] \[\] \[\] \[\]$	\ne	¥	U+2260	relational	normal
$\[\] \[\] \[\] \[\] \[\] \[\] \[\] \[\]$	\nearrow	7	U+2197	relational	normal
\ni∋U+220Brelationalnormal\norm U+2016ordinaryopen/close\nuvU+03BDordinaryoperand\nwarrow\U+2196relationalnormal\odot⊙U+2299binarynormal	\neg	Г	U+00AC	unary	normal
\norm U+2016ordinaryopen/close\nuvU+03BDordinaryoperand\nwarrow\U+2196relationalnormal\odot\U+2299binarynormal	\neq	≠	U+2260	relational	normal
\nuvU+03BDordinaryoperand\nwarrow\U+2196relationalnormal\odot\U+2299binarynormal	\ni	Э	U+220B	relational	normal
\nwarrow\U+2196relationalnormal\odot \odot U+2299binarynormal	\norm		U+2016	ordinary	open/close
\odot①U+2299binarynormal	\nu	ν	U+03BD	ordinary	operand
	\nwarrow	7	U+2196	relational	normal
\of IL+2502 ordinary normal	\odot	\odot	U+2299	binary	normal
	\of		U+2592	ordinary	normal

\oiiint	∰	U+2230	ordinary	nary
\oiint	∯	U+222F	ordinary	nary
\oint	∮	U+222E	ordinary	nary
\Omega	Ω	U+03A9	ordinary	operand
\omega	ω	U+03C9	ordinary	operand
\ominus	θ	U+2296	binary	normal
\open	F	U+251C	ordinary	open
\oplus	\oplus	U+2295	binary	normal
\oslash	\oslash	U+2298	binary	normal
\otimes	\otimes	U+2297	binary	normal
\over	/	U+002F	binarynsp	divide
\overbar	-	U+00AF	ordinary	encl overbar
\overbrace	~	U+23DE	ordinary	stretch over
\overparen	\frown	U+23DC	ordinary	stretch over
\parallel	II	U+2225	relational	normal
\partial	д	U+2202	unary	operand
\phantom	\$	U+27E1	ordinary	encl phantom
\Phi	Φ	U+03A6	ordinary	operand
\phi	φ	U+03D5	ordinary	operand
\Pi	П	U+03A0	ordinary	operand
\pi	π	U+03C0	ordinary	operand
\pm	±	U+00B1	unary/binary	unary/binary
\pppprime	,,,,,	U+2057	ordinary	Unisubsup
\ppprime	""	U+2034	ordinary	Unisubsup
\pprime	"	U+2033	ordinary	Unisubsup
\prcue	≼	U+227C	relational	normal
\prec	\prec	U+227A	relational	normal
\preceq	≤	U+2AAF	relational	normal
\preccurlyeq	≼	U+227C	relational	normal
\prime	'	U+2032	ordinary	Unisubsup
\prod	П	U+220F	ordinary	nary
\propto	∝	U+221D	relational	normal
\Psi	Ψ	U+03A8	ordinary	operand
\psi	ψ	U+03C8	ordinary	operand
\qdrt	$\sqrt[4]{}$	U+221C	open	encl root

\rangle	>	U+27E9	close	close
\ratio	:	U+2236	relational	normal
\rbrace	}	U+007D	close	close
\rbrack]	U+005D	close	close
\rceil	1	U+2309	close	close
\rddots	.:	U+22F0	relational	normal
\Re	R	U+211C	ordinary	operand
\rect		U+25AD	ordinary	encl rect
\rfloor]	U+230B	close	close
\rho	ρ	U+03C1	ordinary	operand
\Rightarrow	↑	U+21D2	relational	stretch horiz
\rightarrow	\rightarrow	U+2192	relational	stretch horiz
\rightharpoondown	7	U+21C1	relational	stretch horiz
\rightharpoonup	-	U+21C0	relational	stretch horiz
\rrect	0	U+25A2	ordinary	encl rnd rect
\sdiv	/	U+2044	binarynsp	divide
\searrow	2	U+2198	relational	normal
\setminus	\	U+2216	binary	normal
\Sigma	Σ	U+03A3	ordinary	operand
\sigma	σ	U+03C3	ordinary	operand
\sim	~	U+223C	relational	normal
\simeq	김	U+2243	relational	normal
\smash	\$	U+2B0D	ordinary	encl phantom
\spadesuit	♠	U+2660	ordinary	normal
\sqcap	П	U+2293	binary	normal
\sqcup		U+2294	binary	normal
\sqrt	\checkmark	U+221A	open	encl root
\sqsubseteq	LI	U+2291	relational	normal
\sqsuperseteq		U+2292	relational	normal
\star	*	U+22C6	binary	normal
\subset	C	U+2282	relational	normal
\subseteq	L	U+2286	relational	normal
\succ	$\boldsymbol{\lambda}$	U+227B	relational	normal
\succeq	≫	U+227D	relational	normal
\sum	Σ	U+2211	ordinary	nary

\superset⊃U+2283relationalnormal\superseteq⊇U+2287relationalnormal\swarrow\superseteqU+2199relationalnormal\tauTU+03C4ordinaryoperand\thereforeU+2234relationalnormal\thereforeU+2234relationalnormal\thereforeU+2038ordinaryoperand\theta0U+0388ordinaryoperand\thickspUU+2005skipnormal\thinspU+2006skipnormal\thinspU+2007binarynspnormal\thinspU+2192relationalstretch horiz\thinspV+2192relationalnormal\thetaU+2192relationalnormal\thetaU+2192relationalnormal\thetaU+2192relationalnormal\thetaU+2192relationalnormal\thetaU+2191ordinaryaccent\underbarU+238Dordinarystretch under\underbarU+2191relationalnormal\underbarU+238Dordinarystretch under\underbarU+238Dordinaryoperand\underbarU+239Dordinarystretch under\underbarU+238Dordinaryoperand<	\ auponact	_	U+2283	rolational	normal
\swarrow\statuU+2199relationalnormal\tau\tau\tau\tauoperand\therefore \therefore U+2234relationalnormal\Theta\theta\thetaU+0398ordinaryoperand\theta\thetaU+2005skipnormal\thickspU+2005skipnormal\thinspU+2006skipnormal\thinspU+2006skipnormal\thinspU+2007binarynspnormal\times \times U+007binarynspnormal\to \rightarrow U+2192relationalstretch horiz\topTU+2244relationalnormal\topTU+2281ordinaryaccent\underbar $_$ U+230Fordinarystretch under\underbarce $_$ U+23DFordinarystretch under\underparen $_$ U+2191relationalnormal\underbarcw \uparrow U+2191relationalnormal\underparen $_$ U+2195relationalnormal\underparen $_$ U+2195relationalnormal\underparen $_$ U+2195relationalnormal\underparen $_$ U+2195relationalnormal\underparen $_$ U+2195relationalnormal\underparen $_$ U+2195relationalnormal\underparen $_$ U+2195relationalnorma					
\tau τ U+03C4ordinaryoperand\therefore \therefore U+2234relationalnormal\Theta Θ U+0398ordinaryoperand\theta θ U+03B8ordinaryoperand\thicksp u U+2005skipnormal\thinsp u U+2006skipnormal\thinsp u U+2006skipnormal\thinsp u U+2006skipnormal\thinsp u U+2007binarynspnormal\times \times U+007binarynspnormal\to \rightarrow U+2192relationalstretch horiz\topTU+2284relationalnormal\topTU+2281ordinaryaccent\underbar $_$ U+230Fordinarystretch under\underbarce $_$ U+23DFordinarystretch under\underparen $_$ U+2101relationalnormal\underparow \uparrow U+2195relationalnormal\underparow \uparrow U+228Ebinaryoperand\underparow \uparrow U+2355ordinaryoperand\underparom $_$ U+3355ordinaryoperand\underparom $□$ U+03C6ordinaryoperand\underparom $□$ U+03C6ordinaryoperand\underparom $□$ U+03C6ordinaryoperand\underparom $□$ U+03C6ordin			-		
\therefore \cdot U+2234relationalnormal\Theta Θ U+0398ordinaryoperand\theta θ U+0388ordinaryoperand\thickspU+2005skipnormal\thinspU+2006skipnormal\thinspU+2006skipnormal\thinspU+2007binarynspnormal\times \times U+0007binarynspnormal\times \times U+0017binarynspnormal\to \rightarrow U+2192relationalstretch horiz\topTU+22A4relationalnormal\tvcc T^* U+20E1ordinaryaccent\underbar $_$ U+23DFordinarystretch under\underbare $_$ U+23DDordinarystretch under\underbarem $_$ U+21D1relationalnormal\updownarrow \uparrow U+21D5relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \downarrow U+2195relationalnormal\updownarrow \downarrow U+2385ordinaryoperand\updownarrow \downarrow U+2195relationalnormal\updownarrow \downarrow U+2195relationalnormal\updownarrow \downarrow U+2195ordinaryoperand\updownarrow \downarrow U+2195ordinaryoperand\updownarrow \downarrow U+0365ordinaryopera	`	7			
\ThetaΘU+0398ordinaryoperand\thetaΘU+03B8ordinaryoperand\thickspU+2005skipnormal\thinspU+2006skipnormal\thinspU+2006skipnormal\tilde~U+0303ordinaryaccent\times×U+0D07binarynspnormal\to→U+2192relationalstretch horiz\topTU+22A4relationalnormal\tvec~U+23DFordinaryaccent\underbar_U+23DFordinarystretch under\underbarce_U+23DFordinarystretch under\underparen_U+21D1relationalnormal\updownarrow↑U+2195relationalnormal\updownarrow↑U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195ordinaryoperand\updownarrow↓U+2195ordinaryoperand\updownarrow↓U+2195ordinaryoperand\updownarrow↓U+2195ordinaryoperand\updownarrow↓U+03C5ordinaryoperand\updownarrow↓U+03C5ordinaryoperand\upd	`	τ		-	-
\theta θ U+03B8ordinaryoperand\thickspU+2005skipnormal\thinspU+2006skipnormal\tilde~U+0007binarynspnormal\tilde~U+0107binarynspnormal\times×U+0007binarynspnormal\to \rightarrow U+2192relationalstretch horiz\topTU+22A4relationalnormal\tvec \leftrightarrow U+20E1ordinaryaccent\underbar_U+23DFordinarystretch under\underbraceU+23DFordinarystretch under\underparenU+21D1relationalnormal\updownarrow \uparrow U+21D1relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \downarrow U+228Ebinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow ξ U+03C5ordinaryoperand\updownarrow \downarrow U+2195relationalnormal\updownarrow \downarrow U+2195ordinaryoperand\updownarrow \downarrow U+2195ordinaryoperand\updownarrow \downarrow U+2195ordinaryo	\therefore	:	U+2234	relational	normal
\thickspUU2005skipnormal\thinspU+2006skipnormal\tilde \sim U+0007binarynspnormal\times \times U+00D7binarynspnormal\to \rightarrow U+2192relationalstretch horiz\topTU+22A4relationalnormal\tvec \leftrightarrow U+20E1ordinaryaccent\underbar_U+23B1ordinaryencl un-\underbraceU+23DFordinarystretch under\upderbarrow \uparrow U+2191relationalnormal\upderbarrow \uparrow U+2191relationalnormal\upderbarrow \uparrow U+2191relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+228Ebinarynormal\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C6ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C5ordinaryoperand\updownarrow \downarrow U+03C6ordinaryoperand\updownarrow ξ	\Theta	Θ	U+0398	ordinary	operand
\thinspU+2006skipnormal\tilde $$ U+0303ordinaryaccent\times \times U+00D7binarynspnormal\to \rightarrow U+2192relationalstretch horiz\topTU+22A4relationalnormal\tvec $\stackrel{\leftrightarrow}{\rightarrow}$ U+20E1ordinaryaccent\underbar_U+23B1ordinaryencl un-\underbraceU+23DFordinarystretch under\underparenU+23DDordinarystretch under\uparrow \uparrow U+2191relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \downarrow U+2385ordinaryoperand\updownarrow ψ U+0366ordinaryoperand\updownarrow </td <td>\theta</td> <td>θ</td> <td>U+03B8</td> <td>ordinary</td> <td>operand</td>	\theta	θ	U+03B8	ordinary	operand
\tilde $\tilde{~}$ U+0303ordinaryaccent\times \times U+00D7binarynspnormal\to \rightarrow U+2192relationalstretch horiz\topTU+22A4relationalnormal\tvec $\stackrel{\leftrightarrow}{}$ U+20E1ordinaryaccent\underbar $_$ U+22B1ordinaryencl un-\underbrace \bigcup U+23DFordinarystretch under\underparen \bigcirc U+23DDordinarystretch under\uparrow \uparrow U+21D1relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+23B5ordinaryoperand\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \uparrow U+2385ordinaryoperand\updownarrow \downarrow U+0355ordinaryoperand\updownarrow \downarrow U+0355ordinaryoperand\updownarrow ς U+0355ordinaryoperand\updownarrow ς U+0356ordinaryoperand\updownarrow ς U+0356ordinaryoperand\updownarrow ς U+0356ordinaryoperand\updownarrow ς U+0356ordinary<	\thicksp		U+2005	skip	normal
Allace $0+0503$ ordinaryaccent\times \times $U+00D7$ binarynspnormal\to \rightarrow $U+2192$ relationalstretch horiz\topT $U+22A4$ relationalnormal\tvec $\stackrel{\leftarrow}{\rightarrow}$ $U+20E1$ ordinaryaccent\underbar $_$ $U+2581$ ordinaryencl un-\underbarce $_$ $U+23DF$ ordinarystretch under\underparen $_$ $U+23DD$ ordinarystretch under\updownarrow \Uparrow $U+21D1$ relationalnormal\updownarrow \Uparrow $U+2195$ relationalnormal\updownarrow \Uparrow $U+2195$ relationalnormal\updownarrow \Uparrow $U+2195$ relationalnormal\updownarrow \Uparrow $U+2195$ relationalnormal\updownarrow \updownarrow $U+2195$ relationalnormal\updownarrow \checkmark $U+2385$ ordinaryoperand\updownarrow \checkmark $U+0365$ ordinaryoperand\updownarrow \pounds $U+0326$ ordinaryoperand\updownaro \blacklozenge <	\thinsp		U+2006	skip	normal
\to→U+2192relationalstretch horiz\topTU+2192relationalnormal\tvec $\stackrel{\leftrightarrow}{}$ U+20E1ordinaryaccent\underbar_U+23B1ordinaryencl un-\underparenU+23DFordinarystretch under\uparrow \uparrow U+21D1relationalnormal\upderbarceU+23DDordinarystretch under\uparrow \uparrow U+21D1relationalnormal\updownarrow \uparrow U+2195relationalnormal\updownarrow \downarrow U+2195relationalnormal\updownarrow \downarrow U+228Ebinaryoperand\upsilon Υ U+03A5ordinaryoperand\varepsilon ε U+03B5ordinaryoperand\varphi ϕ U+03C6ordinaryoperand\varphi φ U+03D6ordinaryoperand\varphi ϕ U+03D6ordinaryoperand\varphi ϕ U+03D6ordinaryoperand\varphi φ U+03D6ordinaryoperand\varphi φ U+03D1ordinaryoperand\varphi φ U+03D2ordinaryoperand\varphi φ U+03D6ordinaryoperand\varphi φ U+03D6ordinaryoperand\varphi φ U+03D1ordinaryoperand\varphi φ U+03D2ordinary<	\tilde	~	U+0303	ordinary	accent
\topTU+22A4relationalnormal\tvec→U+20E1ordinaryaccent\underbar_U+2581ordinaryencl un-\underbraceU+23DFordinarystretch under\underparenU+23DDordinarystretch under\uparrow↑U+21D1relationalnormal\updownarrow↑U+2191relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195ordinaryoperand\updownarrow↓U+2385ordinaryoperand\updownarrow↓U+03C5ordinaryoperand\updownarrow€U+03C6ordinaryoperand\updownarepsilon€U+03C6ordinaryoperand\varpi∅U+03C1ordinaryoperand\varpi∅U+03C2ordinaryoperand\varpi∅U+03C2ordinaryoperand\varpi∅U+03C2ordinaryoperand\varpi∅U+03C1ordina	\times	×	U+00D7	binarynsp	normal
\tvec\thereforeU+20E1ordinaryaccent\underbarU+23DFordinaryencl un-\underbraceU+23DFordinarystretch under\underparenU+23DDordinarystretch under\Uparrow\thetaU+21D1relationalnormal\updownarrow\thetaU+21D5relationalnormal\updownarrow\thetaU+2195relationalnormal\updownarrow\thetaU+2185relationalnormal\updownarrow\thetaU+2185relationalnormal\updownarrow\thetaU+2185relationalnormal\updownarrow\thetaU+2185relationalnormal\updownarrow\thetaU+2185ordinaryoperand\updownarrow\thetaU+2185ordinaryoperand\updownarrow\thetaU+2185ordinaryoperand\updownarrow\thetaU+2185ordinaryoperand\updownarrow\thetaU+0365ordinaryoperand\updownarrow\thetaU+0385ordinaryoperand\updownarrow\thetaU+0385ordinaryoperand\updownarrow\thetaU+0386ordinaryoperand\updownarrow\thetaU+0386ordinaryoperand\updownarrow\thetaU+0385ordinaryoperand\updownarrow\thetaU+0386ordinaryoperand <td< td=""><td>\to</td><td>\rightarrow</td><td>U+2192</td><td>relational</td><td>stretch horiz</td></td<>	\to	\rightarrow	U+2192	relational	stretch horiz
\tvec0+20E1ordinaryaccent\underbar_U+2581ordinaryencl un-\underbraceU+23DFordinarystretch under\underparenU+23DDordinarystretch under\uparrow↑U+21D1relationalnormal\updownarrow↑U+2191relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+218Ebinarynormal\updownarrow↓U+238Eordinaryoperand\updownarrow↓U+03A5ordinaryoperand\updownarrow↓U+03A5ordinaryoperand\updownarrow↓U+03C6ordinaryoperand\updownarrow↓U+03C6ordinaryoperand\updownarrow↓U+03C6ordinaryoperand\updownarrow↓U+03C6ordinaryoperand\updownarrow↓U+03C6ordinaryoperand\updownarrow↓U+03D6ordinaryoperand\updownarrow↓U+03C6ordinaryoperand\updownarrow↓U+03D6ordinaryoperand\updownarrow↓U+03D1ordinaryoperand\updownarrow↓U+03D2ordinaryoperand\updownarrow↓U+03D2	\top	Т	U+22A4	relational	normal
\underbraceU+23DFordinarystretch under\underparenU+23DDordinarystretch under\Uparrow↑U+21D1relationalnormal\uparrow↑U+21P1relationalnormal\updownarrow↑U+21D5relationalnormal\updownarrow↑U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5relationalnormal\updownarrow↓U+21P5ordinaryoperand\updownarrow↓U+03C5ordinaryoperand\updownarrow₽U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6	\tvec	\leftrightarrow	U+20E1	ordinary	accent
\underparenU+23DDordinarystretch under\Uparrow↑U+23DDrelationalnormal\uparrow↑U+21D1relationalnormal\Updownarrow\$U+21D5relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+03C5ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03F1ordinaryoperand\varphi\$U+03C1ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C1ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2<	\underbar	_	U+2581	ordinary	encl un-
\underparenU+23DDordinarystretch under\Uparrow↑U+21D1relationalnormal\uparrow↑U+2191relationalnormal\Updownarrow↑U+21D5relationalnormal\updownarrow↑U+2195relationalnormal\updownarrow↑U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195relationalnormal\updownarrow↓U+2195ordinaryoperand\updownarrow↓U+03A5ordinaryoperand\upsilon↓U+03C5ordinaryoperand\upsilon↓U+03C5ordinaryoperand\varepsilon€U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$\$U+03C6ordinaryop	\underbrace	<u></u>	U+23DF	ordinary	stretch under
\uparrow↑U+2191relationalnormal\Updownarrow\$U+21D5relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195ordinaryoperand\updownarrow\$U+03C5ordinaryoperand\updownarrow\$U+03C5ordinaryoperand\updownarrow\$U+03C6ordinaryoperand\varepsilon\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03D6ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+03C2ordinaryoperand\varphi\$U+2502ordinary	\underparen		U+23DD	ordinary	stretch under
\Updownarrow\$U+21D5relationalnormal\updownarrow\$U+2195relationalnormal\updownarrow\$U+2195relationalnormal\updus\$U+228Ebinarynormal\upsilonYU+03A5ordinaryoperand\upsilon\$U+03C5ordinaryoperand\varepsilon\$U+03C5ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03C6ordinaryoperand\varphi\$U+03D6ordinaryoperand\varphi\$U+03D1ordinaryoperand\varho\$U+03D1ordinaryoperand\varheta\$U+03D1ordinaryoperand\varheta\$U+03D1ordinaryoperand\varheta\$U+03D1ordinaryoperand\varheta\$U+03D1ordinaryoperand\varheta\$U+03D1ordinaryoperand\varheta\$U+03D1ordinaryoperand\varheta\$U+22E2relationalstretch horz\vdots\$U+22D7ordinaryaccent	\Uparrow	Î	U+21D1	relational	normal
\updownarrow↓U+2195relationalnormal\uplus⊌U+228Ebinarynormal\UpsilonYU+03A5ordinaryoperand\upsilonvU+03C5ordinaryoperand\varepsilonεU+03B5ordinaryoperand\varphiφU+03C6ordinaryoperand\varphiφU+03C6ordinaryoperand\varphiφU+03C6ordinaryoperand\varphiφU+03C6ordinaryoperand\varphiφU+03D6ordinaryoperand\varphiφU+03D1ordinaryoperand\varthetaθU+03D1ordinaryoperand\vbar U+2502ordinaryoperand\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec-U+20D7ordinaryaccent	\uparrow	↑	U+2191	relational	normal
\uplus⊌U+228Ebinarynormal\UpsilonΥU+03A5ordinaryoperand\upsilonυU+03C5ordinaryoperand\varepsilonεU+03B5ordinaryoperand\varphiφU+03C6ordinaryoperand\varphiφU+03C6ordinaryoperand\varpi∞U+03C6ordinaryoperand\varpi∞U+03D6ordinaryoperand\varsigmaςU+03D1ordinaryoperand\varthetaθU+03D1ordinaryoperand\vbar U+2502ordinaryoperand\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec-U+20D7ordinaryaccent	\Updownarrow	\$	U+21D5	relational	normal
\UpsilonYU+03A5ordinaryoperand\upsilonυU+03C5ordinaryoperand\varepsilonεU+03B5ordinaryoperand\varphiφU+03C6ordinaryoperand\varphiφU+03C6ordinaryoperand\varpi\overline\overlineordinaryoperand\varpi\overlineU+03D6ordinaryoperand\varsigma\squteU+03F1ordinaryoperand\varsigma\squteU+03C2ordinaryoperand\vartheta\overlineU+03D1ordinaryoperand\varsigma\squteU+03D1ordinaryoperand\varsigma\squteU+03D1ordinaryoperand\varsigma\squteU+03D2ordinaryoperand\varsigma\squteU+03D1ordinaryoperand\varsigma\squteU+03D1ordinaryoperand\varsigma\squteU+22A2relationalstretch horz\vdots\varsigmaU+22EErelationalnormal\vec-\fractorU+20D7ordinaryaccent	\updownarrow	\$	U+2195	relational	normal
\upsilonυU+03C5ordinaryoperand\varepsilonεU+03B5ordinaryoperand\varphiφU+03C6ordinaryoperand\varpi\oversigmaU+03D6ordinaryoperand\varrhoQU+03F1ordinaryoperand\varsigma\sigmaU+03C2ordinaryoperand\vartheta\varthetaU+03D1ordinaryoperand\vbar U+2502ordinaryoperand\vdash\-U+22A2relationalstretch horz\vdots!U+22EErelationalnormal\vec-U+20D7ordinaryaccent	\uplus	Ŀ	U+228E	binary	normal
\varepsilonεU+03B5ordinaryoperand\varphiφU+03C6ordinaryoperand\varpi\varthoQU+03D6ordinaryoperand\varrhoQU+03F1ordinaryoperand\varsigma\sigmaU+03C2ordinaryoperand\vartheta\vartheta\vartheta\varthetaordinaryoperand\vbar U+03D1ordinaryoperand\vbar U+2502ordinaryoperand\vdash\-U+22A2relationalstretch horz\vdots:U+22EErelationalnormal\vec-U+20D7ordinaryaccent	\Upsilon	Ŷ	U+03A5	ordinary	operand
\varphiφU+03C6ordinaryoperand\varpiωU+03D6ordinaryoperand\varrhoQU+03F1ordinaryoperand\varsigmaςU+03C2ordinaryoperand\varthetaϑU+03D1ordinaryoperand\vbarIU+2502ordinaryist delims\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vecIU+20D7ordinaryaccent	\upsilon	υ	U+03C5	ordinary	operand
\varpi\varpi\varpi\varpi\varpi\varpi\varrho\varpi\varpiordinaryoperand\varsigma\sigma\sigma\varpiordinaryoperand\varsigma\sigma\sigma\varpiordinaryoperand\vartheta\varpiU+03D1ordinaryoperand\vbar U+2502ordinarylist delims\vdash\-U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec''U+20D7ordinaryaccent	\varepsilon	3	U+03B5	ordinary	operand
\varrhoQU+03F1ordinaryoperand\varsigma\$U+03C2ordinaryoperand\vartheta\$U+03D1ordinaryoperand\vbar U+2502ordinarylist delims\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec✓U+20D7ordinaryaccent	\varphi	φ	U+03C6	ordinary	operand
\varsigmaçU+03C2ordinaryoperand\varthetaϑU+03D1ordinaryoperand\vbarIU+2502ordinarylist delims\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec✓U+20D7ordinaryaccent	\varpi	ω	U+03D6	ordinary	operand
\varthetaθU+03D1ordinaryoperand\vbar U+2502ordinarylist delims\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec⁻U+20D7ordinaryaccent	\varrho	Q	U+03F1	ordinary	operand
\varthetaϑU+03D1ordinaryoperand\vbar□U+2502ordinarylist delims\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec⁻U+20D7ordinaryaccent	\varsigma	ς	U+03C2	ordinary	operand
\vdash⊢U+22A2relationalstretch horz\vdots⋮U+22EErelationalnormal\vec`U+20D7ordinaryaccent	\vartheta		U+03D1	ordinary	operand
\vdots:U+22EErelationalnormal\vecU+20D7ordinaryaccent	\vbar		U+2502	ordinary	list delims
\vecImage: U+20D7ordinaryaccent	\vdash	F	U+22A2	relational	stretch horz
Vec 0+20D7 ordinary accent	\vdots	:	U+22EE	relational	normal
	\vec	→	U+20D7	ordinary	accent
		V	U+2228	binary	normal

\Vert		U+2016	ordinary	open/close
\vert		U+007C	ordinary	open/close
\vphantom	Û	U+21F3	relational	encl phantom
\vthicksp		U+2004	skip	normal
\wedge	Λ	U+2227	binary	normal
\wp	80	U+2118	ordinary	operand
\wr	2	U+2240	binary	normal
\Xi	[1]	U+039E	ordinary	operand
\xi	ξ	U+03BE	ordinary	operand
\zeta	ζ	U+03B6	ordinary	operand
\zwnj		U+200C	ordinary	normal
\zwsp		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 Version 2 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 expressions as their syntax becomes unambiguous.

In this document, features added in Version 3 are identified as such. These features are mostly 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 nth 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 a number of topics up to date and using the name UnicodeMath instead of Unicode linear format.

References

- 1. The Unicode Standard http://www.unicode.org/versions/latest/; see also².
- Barbara Beeton, Asmus Freytag, Murray Sargent III, Unicode Technical Report #25 "Unicode Support for Mathematics", <u>http://www.unicode.org/reports/tr25</u>
- 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) <u>http://www.w3.org/Math/</u>
- 6. For example, UnicodeMath is used for keyboard entry of mathematical expressions in Microsoft Word, PowerPoint, OneNote and Excel.
- Bertrand Russell, in his Introduction to *Tractatus Logico-Philosophicus* by Lugwig Wittgenstein, Routledge and Kegan Paul, London 1922 (also currently available at <u>http://www.kfs.org/~jonathan/witt/tlph.html</u>).
- 8. PS Technical Word Processor, Scroll Systems, Inc. (1989). This WP used a non-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).
- 11. Alexander Mamishev, Murray Sargent (2103), *Creating Research and Scientific Documents with Microsoft Word*, Microsoft Press, Redmond, WA.

This document was prepared using Microsoft Word 2016 with Cambria and Cambria Math fonts.