

Modul:Unicode data/blocks

Ausgabe: 27.06.2026

Letzte Änderung: 23.02.2022

Seite von

Inhaltsverzeichnis

- [1. Modul:Unicode data/blocks](#)
- [2. Modul:Unicode data](#)

Modul:Unicode data/blocks

Insert non-formatted text here Contains data on [Unicode blocks](#) for [Module:Unicode data](#) derived from [Blocks.txt](#) in the Unicode Character Database.

```
-- Compiled from http://www.unicode.org/Public/UNIDATA/Blocks.txt.  
local blocks = {
```

```
{ 0x000000, 0x00007F, "Basic Latin" },  
{ 0x000080, 0x0000FF, "Latin-1 Supplement" },  
{ 0x000100, 0x00017F, "Latin Extended-A" },  
{ 0x000180, 0x00024F, "Latin Extended-B" },  
{ 0x000250, 0x0002AF, "IPA Extensions" },  
{ 0x0002B0, 0x0002FF, "Spacing Modifier Letters" },  
{ 0x000300, 0x00036F, "Combining Diacritical Marks" },  
{ 0x000370, 0x0003FF, "Greek and Coptic" },  
{ 0x000400, 0x0004FF, "Cyrillic" },  
{ 0x000500, 0x00052F, "Cyrillic Supplement" },  
{ 0x000530, 0x00058F, "Armenian" },  
{ 0x000590, 0x0005FF, "Hebrew" },  
{ 0x000600, 0x0006FF, "Arabic" },  
{ 0x000700, 0x00074F, "Syriac" },  
{ 0x000750, 0x00077F, "Arabic Supplement" },  
{ 0x000780, 0x0007BF, "Thaana" },  
{ 0x0007C0, 0x0007FF, "Nko" },  
{ 0x000800, 0x00083F, "Samaritan" },  
{ 0x000840, 0x00085F, "Mandaic" },  
{ 0x000860, 0x00086F, "Syriac Supplement" },  
{ 0x000870, 0x00089F, "Arabic Extended-B" },  
{ 0x0008A0, 0x0008FF, "Arabic Extended-A" },  
{ 0x000900, 0x00097F, "Devanagari" },  
{ 0x000980, 0x0009FF, "Bengali" },  
{ 0x000A00, 0x000A7F, "Gurmukhi" },  
{ 0x000A80, 0x000AFF, "Gujarati" },  
{ 0x000B00, 0x000B7F, "Oriya" },  
{ 0x000B80, 0x000BFF, "Tamil" },  
{ 0x000C00, 0x000C7F, "Telugu" },  
{ 0x000C80, 0x000CFF, "Kannada" },  
{ 0x000D00, 0x000D7F, "Malayalam" },  
{ 0x000D80, 0x000DFF, "Sinhala" },  
{ 0x000E00, 0x000E7F, "Thai" },  
{ 0x000E80, 0x000EFF, "Lao" },  
{ 0x000F00, 0x000FFF, "Tibetan" },
```

0x001000	0x00109F	"Myanmar"
0x0010A0	0x0010FF	"Georgian"
0x001100	0x0011FF	"Hangul Jamo"
0x001200	0x00137F	"Ethiopic"
0x001380	0x00139F	"Ethiopic Supplement"
0x0013A0	0x0013FF	"Cherokee"
0x001400	0x00167F	"Unified Canadian Aboriginal Syllabics"
0x001680	0x00169F	"Ogham"
0x0016A0	0x0016FF	"Runic"
0x001700	0x00171F	"Tagalog"
0x001720	0x00173F	"Hanunoo"
0x001740	0x00175F	"Buhid"
0x001760	0x00177F	"Tagbanwa"
0x001780	0x0017FF	"Khmer"
0x001800	0x0018AF	"Mongolian"
0x0018B0	0x0018FF	"Unified Canadian Aboriginal Syllabics Extended"
0x001900	0x00194F	"Limbu"
0x001950	0x00197F	"Tai Le"
0x001980	0x0019DF	"New Tai Lue"
0x0019E0	0x0019FF	"Khmer Symbols"
0x001A00	0x001A1F	"Buginese"
0x001A20	0x001AAF	"Tai Tham"
0x001AB0	0x001AFF	"Combining Diacritical Marks Extended"
0x001B00	0x001B7F	"Balinese"
0x001B80	0x001BBF	"Sundanese"
0x001BC0	0x001BFF	"Batak"
0x001C00	0x001C4F	"Lepcha"
0x001C50	0x001C7F	"Ol Chiki"
0x001C80	0x001C8F	"Cyrillic Extended-C"
0x001C90	0x001CBF	"Georgian Extended"
0x001CC0	0x001CCF	"Sundanese Supplement"
0x001CD0	0x001CFF	"Vedic Extensions"
0x001D00	0x001D7F	"Phonetic Extensions"
0x001D80	0x001DBF	"Phonetic Extensions Supplement"
0x001DC0	0x001DFF	"Combining Diacritical Marks Supplement"
0x001E00	0x001EFF	"Latin Extended Additional"
0x001F00	0x001FFF	"Greek Extended"
0x002000	0x00206F	"General Punctuation"
0x002070	0x00209F	"Superscripts and Subscripts"
0x0020A0	0x0020CF	"Currency Symbols"
0x0020D0	0x0020FF	"Combining Diacritical Marks for Symbols"
0x002100	0x00214F	"Letterlike Symbols"
0x002150	0x00218F	"Number Forms"
0x002190	0x0021FF	"Arrows"
0x002200	0x0022FF	"Mathematical Operators"
0x002300	0x0023FF	"Miscellaneous Technical"
0x002400	0x00243F	"Control Pictures"
0x002440	0x00245F	"Optical Character Recognition"
0x002460	0x0024FF	"Enclosed Alphanumerics"
0x002500	0x00257F	"Box Drawing"
0x002580	0x00259F	"Block Elements"
0x0025A0	0x0025FF	"Geometric Shapes"
0x002600	0x0026FF	"Miscellaneous Symbols"
0x002700	0x0027BF	"Dingbats"
0x0027C0	0x0027EF	"Miscellaneous Mathematical Symbols-A"
0x0027F0	0x0027FF	"Supplemental Arrows-A"
0x002800	0x0028FF	"Braille Patterns"
0x002900	0x00297F	"Supplemental Arrows-B"
0x002980	0x0029FF	"Miscellaneous Mathematical Symbols-B"
0x002A00	0x002AFF	"Supplemental Mathematical Operators"
0x002B00	0x002BFF	"Miscellaneous Symbols and Arrows"
0x002C00	0x002C5F	"Glagolitic"
0x002C60	0x002C7F	"Latin Extended-C"
0x002C80	0x002CFF	"Coptic"
0x002D00	0x002D2F	"Georgian Supplement"
0x002D30	0x002D7F	"Tifinagh"
0x002D80	0x002DDF	"Ethiopic Extended"

0x002DE0	0x002DFF	"Cyrillic Extended-A"
0x002E00	0x002E7F	"Supplemental Punctuation"
0x002E80	0x002EFF	"CJK Radicals Supplement"
0x002F00	0x002FDF	"Kangxi Radicals"
0x002FF0	0x002FFF	"Ideographic Description Characters"
0x003000	0x00303F	"CJK Symbols and Punctuation"
0x003040	0x00309F	"Hiragana"
0x0030A0	0x0030FF	"Katakana"
0x003100	0x00312F	"Bopomofo"
0x003130	0x00318F	"Hangul Compatibility Jamo"
0x003190	0x00319F	"Kanbun"
0x0031A0	0x0031BF	"Bopomofo Extended"
0x0031C0	0x0031EF	"CJK Strokes"
0x0031F0	0x0031FF	"Katakana Phonetic Extensions"
0x003200	0x0032FF	"Enclosed CJK Letters and Months"
0x003300	0x0033FF	"CJK Compatibility"
0x003400	0x004DBF	"CJK Unified Ideographs Extension A"
0x004DC0	0x004DFF	"Yijing Hexagram Symbols"
0x004E00	0x009FFF	"CJK Unified Ideographs"
0x00A000	0x00A48F	"Yi Syllables"
0x00A490	0x00A4CF	"Yi Radicals"
0x00A4D0	0x00A4FF	"Lisu"
0x00A500	0x00A63F	"Vai"
0x00A640	0x00A69F	"Cyrillic Extended-B"
0x00A6A0	0x00A6FF	"Bamum"
0x00A700	0x00A71F	"Modifier Tone Letters"
0x00A720	0x00A7FF	"Latin Extended-D"
0x00A800	0x00A82F	"Syloti Nagri"
0x00A830	0x00A83F	"Common Indic Number Forms"
0x00A840	0x00A87F	"Phags-pa"
0x00A880	0x00A8DF	"Saurashtra"
0x00A8E0	0x00A8FF	"Devanagari Extended"
0x00A900	0x00A92F	"Kayah Li"
0x00A930	0x00A95F	"Rejang"
0x00A960	0x00A97F	"Hangul Jamo Extended-A"
0x00A980	0x00A9DF	"Javanese"
0x00A9E0	0x00A9FF	"Myanmar Extended-B"
0x00AA00	0x00AA5F	"Cham"
0x00AA60	0x00AA7F	"Myanmar Extended-A"
0x00AA80	0x00AADF	"Tai Viet"
0x00AAE0	0x00AAFF	"Meetei Mayek Extensions"
0x00AB00	0x00AB2F	"Ethiopic Extended-A"
0x00AB30	0x00AB6F	"Latin Extended-E"
0x00AB70	0x00ABBF	"Cherokee Supplement"
0x00ABC0	0x00ABFF	"Meetei Mayek"
0x00AC00	0x00D7AF	"Hangul Syllables"
0x00D7B0	0x00D7FF	"Hangul Jamo Extended-B"
0x00D800	0x00DB7F	"High Surrogates"
0x00DB80	0x00DBFF	"High Private Use Surrogates"
0x00DC00	0x00DFFF	"Low Surrogates"
0x00E000	0x00F8FF	"Private Use Area"
0x00F900	0x00FAFF	"CJK Compatibility Ideographs"
0x00FB00	0x00FB4F	"Alphabetic Presentation Forms"
0x00FB50	0x00FDFD	"Arabic Presentation Forms-A"
0x00FE00	0x00FE0F	"Variation Selectors"
0x00FE10	0x00FE1F	"Vertical Forms"
0x00FE20	0x00FE2F	"Combining Half Marks"
0x00FE30	0x00FE4F	"CJK Compatibility Forms"
0x00FE50	0x00FE6F	"Small Form Variants"
0x00FE70	0x00FEFF	"Arabic Presentation Forms-B"
0x00FF00	0x00FFEF	"Halfwidth and Fullwidth Forms"
0x00FFF0	0x00FFFF	"Specials"
0x010000	0x01007F	"Linear B Syllabary"
0x010080	0x0100FF	"Linear B Ideograms"
0x010100	0x01013F	"Aegean Numbers"
0x010140	0x01018F	"Ancient Greek Numbers"
0x010190	0x0101CF	"Ancient Symbols"

0x0101D0	0x0101FF	"Phaistos Disc"
0x010280	0x01029F	"Lycian"
0x0102A0	0x0102DF	"Carian"
0x0102E0	0x0102FF	"Coptic Epact Numbers"
0x010300	0x01032F	"Old Italic"
0x010330	0x01034F	"Gothic"
0x010350	0x01037F	"Old Permic"
0x010380	0x01039F	"Ugaritic"
0x0103A0	0x0103DF	"Old Persian"
0x010400	0x01044F	"Deseret"
0x010450	0x01047F	"Shavian"
0x010480	0x0104AF	"Osmanya"
0x0104B0	0x0104FF	"Osage"
0x010500	0x01052F	"Elbasan"
0x010530	0x01056F	"Caucasian Albanian"
0x010570	0x0105BF	"Vithkuqi"
0x010600	0x01077F	"Linear A"
0x010780	0x0107BF	"Latin Extended-F"
0x010800	0x01083F	"Cypriot Syllabary"
0x010840	0x01085F	"Imperial Aramaic"
0x010860	0x01087F	"Palmyrene"
0x010880	0x0108AF	"Nabataean"
0x0108E0	0x0108FF	"Hatran"
0x010900	0x01091F	"Phoenician"
0x010920	0x01093F	"Lydian"
0x010980	0x01099F	"Meroitic Hieroglyphs"
0x0109A0	0x0109FF	"Meroitic Cursive"
0x010A00	0x010A5F	"Kharoshthi"
0x010A60	0x010A7F	"Old South Arabian"
0x010A80	0x010A9F	"Old North Arabian"
0x010AC0	0x010AFF	"Manichaean"
0x010B00	0x010B3F	"Avestan"
0x010B40	0x010B5F	"Inscriptional Parthian"
0x010B60	0x010B7F	"Inscriptional Pahlavi"
0x010B80	0x010BAF	"Psalter Pahlavi"
0x010C00	0x010C4F	"Old Turkic"
0x010C80	0x010CFF	"Old Hungarian"
0x010D00	0x010D3F	"Hanifi Rohingya"
0x010E60	0x010E7F	"Rumi Numeral Symbols"
0x010E80	0x010EBF	"Yezidi"
0x010F00	0x010F2F	"Old Sogdian"
0x010F30	0x010F6F	"Sogdian"
0x010F70	0x010FAF	"Old Uyghur"
0x010FB0	0x010FDF	"Chorasmian"
0x010FE0	0x010FFF	"Elymaic"
0x011000	0x01107F	"Brahmi"
0x011080	0x0110CF	"Kaithi"
0x0110D0	0x0110FF	"Sora Sompeng"
0x011100	0x01114F	"Chakma"
0x011150	0x01117F	"Mahajani"
0x011180	0x0111DF	"Sharada"
0x0111E0	0x0111FF	"Sinhala Archaic Numbers"
0x011200	0x01124F	"Khojki"
0x011280	0x0112AF	"Multani"
0x0112B0	0x0112FF	"Khudawadi"
0x011300	0x01137F	"Grantha"
0x011400	0x01147F	"Newa"
0x011480	0x0114DF	"Tirhuta"
0x011580	0x0115FF	"Siddham"
0x011600	0x01165F	"Modi"
0x011660	0x01167F	"Mongolian Supplement"
0x011680	0x0116CF	"Takri"
0x011700	0x01174F	"Ahom"
0x011800	0x01184F	"Dogra"
0x0118A0	0x0118FF	"Warang Citi"
0x011900	0x01195F	"Dives Akuru"
0x0119A0	0x0119FF	"Nandinagari"

{ 0x011A00, 0x011A4F, "Zanabazar Square" }	}
{ 0x011A50, 0x011AAF, "Soyombo" }	}
{ 0x011AB0, 0x011ABF, "Unified Canadian Aboriginal Syllabics Extended-A" }	}
{ 0x011AC0, 0x011AFF, "Pau Cin Hau" }	}
{ 0x011C00, 0x011C6F, "Bhaiksuki" }	}
{ 0x011C70, 0x011CBF, "Marchen" }	}
{ 0x011D00, 0x011D5F, "Masaram Gondi" }	}
{ 0x011D60, 0x011DAF, "Gunjala Gondi" }	}
{ 0x011EE0, 0x011EFF, "Makasar" }	}
{ 0x011FB0, 0x011FBF, "Lisu Supplement" }	}
{ 0x011FC0, 0x011FFF, "Tamil Supplement" }	}
{ 0x012000, 0x0123FF, "Cuneiform" }	}
{ 0x012400, 0x01247F, "Cuneiform Numbers and Punctuation" }	}
{ 0x012480, 0x01254F, "Early Dynastic Cuneiform" }	}
{ 0x012F90, 0x012FFF, "Cypro-Minoan" }	}
{ 0x013000, 0x01342F, "Egyptian Hieroglyphs" }	}
{ 0x013430, 0x01343F, "Egyptian Hieroglyph Format Controls" }	}
{ 0x014400, 0x01467F, "Anatolian Hieroglyphs" }	}
{ 0x016800, 0x016A3F, "Bamum Supplement" }	}
{ 0x016A40, 0x016A6F, "Mro" }	}
{ 0x016A70, 0x016ACF, "Tangsa" }	}
{ 0x016AD0, 0x016AFF, "Bassa Vah" }	}
{ 0x016B00, 0x016B8F, "Pahawh Hmong" }	}
{ 0x016E40, 0x016E9F, "Medefaidrin" }	}
{ 0x016F00, 0x016F9F, "Miao" }	}
{ 0x016FE0, 0x016FFF, "Ideographic Symbols and Punctuation" }	}
{ 0x017000, 0x0187FF, "Tangut" }	}
{ 0x018800, 0x018AFF, "Tangut Components" }	}
{ 0x018B00, 0x018CFF, "Khitan Small Script" }	}
{ 0x018D00, 0x018D7F, "Tangut Supplement" }	}
{ 0x01AFF0, 0x01AFFF, "Kana Extended-B" }	}
{ 0x01B000, 0x01B0FF, "Kana Supplement" }	}
{ 0x01B100, 0x01B12F, "Kana Extended-A" }	}
{ 0x01B130, 0x01B16F, "Small Kana Extension" }	}
{ 0x01B170, 0x01B2FF, "Nushu" }	}
{ 0x01BC00, 0x01BC9F, "Duployan" }	}
{ 0x01BCA0, 0x01BCAF, "Shorthand Format Controls" }	}
{ 0x01CF00, 0x01CFCF, "Znamenny Musical Notation" }	}
{ 0x01D000, 0x01D0FF, "Byzantine Musical Symbols" }	}
{ 0x01D100, 0x01D1FF, "Musical Symbols" }	}
{ 0x01D200, 0x01D24F, "Ancient Greek Musical Notation" }	}
{ 0x01D2E0, 0x01D2FF, "Mayan Numerals" }	}
{ 0x01D300, 0x01D35F, "Tai Xuan Jing Symbols" }	}
{ 0x01D360, 0x01D37F, "Counting Rod Numerals" }	}
{ 0x01D400, 0x01D7FF, "Mathematical Alphanumeric Symbols" }	}
{ 0x01D800, 0x01DAAF, "Sutton SignWriting" }	}
{ 0x01DF00, 0x01DFFF, "Latin Extended-G" }	}
{ 0x01E000, 0x01E02F, "Glagolitic Supplement" }	}
{ 0x01E100, 0x01E14F, "Nyiakeng Puachue Hmong" }	}
{ 0x01E290, 0x01E2BF, "Toto" }	}
{ 0x01E2C0, 0x01E2FF, "Wancho" }	}
{ 0x01E7E0, 0x01E7FF, "Ethiopic Extended-B" }	}
{ 0x01E800, 0x01E8DF, "Mende Kikakui" }	}
{ 0x01E900, 0x01E95F, "Adlam" }	}
{ 0x01EC70, 0x01ECBF, "Indic Siyaq Numbers" }	}
{ 0x01ED00, 0x01ED4F, "Ottoman Siyaq Numbers" }	}
{ 0x01EE00, 0x01EEFF, "Arabic Mathematical Alphabetic Symbols" }	}
{ 0x01F000, 0x01F02F, "Mahjong Tiles" }	}
{ 0x01F030, 0x01F09F, "Domino Tiles" }	}
{ 0x01F0A0, 0x01F0FF, "Playing Cards" }	}
{ 0x01F100, 0x01F1FF, "Enclosed Alphanumeric Supplement" }	}
{ 0x01F200, 0x01F2FF, "Enclosed Ideographic Supplement" }	}
{ 0x01F300, 0x01F5FF, "Miscellaneous Symbols and Pictographs" }	}
{ 0x01F600, 0x01F64F, "Emoticons" }	}
{ 0x01F650, 0x01F67F, "Ornamental Dingbats" }	}
{ 0x01F680, 0x01F6FF, "Transport and Map Symbols" }	}
{ 0x01F700, 0x01F77F, "Alchemical Symbols" }	}

```

    { 0x01F780, 0x01F7FF, "Geometric Shapes Extended" },
    { 0x01F800, 0x01F8FF, "Supplemental Arrows-C" },
    { 0x01F900, 0x01F9FF, "Supplemental Symbols and Pictographs" },
    { 0x01FA00, 0x01FA6F, "Chess Symbols" },
    { 0x01FA70, 0x01FAFF, "Symbols and Pictographs Extended-A" },
    { 0x01FB00, 0x01FBFF, "Symbols for Legacy Computing" },
    { 0x020000, 0x02A6DF, "CJK Unified Ideographs Extension B" },
    { 0x02A700, 0x02B73F, "CJK Unified Ideographs Extension C" },
    { 0x02B740, 0x02B81F, "CJK Unified Ideographs Extension D" },
    { 0x02B820, 0x02CEAF, "CJK Unified Ideographs Extension E" },
    { 0x02CEB0, 0x02EBEF, "CJK Unified Ideographs Extension F" },
    { 0x02F800, 0x02FA1F, "CJK Compatibility Ideographs Supplement" },
    { 0x030000, 0x03134F, "CJK Unified Ideographs Extension G" },
    { 0x0E0000, 0x0E007F, "Tags" },
    { 0x0E0100, 0x0E01EF, "Variation Selectors Supplement" },
    { 0x0F0000, 0x0FFFFFFF, "Supplementary Private Use Area-A" },
    { 0x100000, 0x10FFFF, "Supplementary Private Use Area-B" },
}
blocks.length = #blocks

return blocks

```

Modul:Unicode data

Inhaltsverzeichnis

- [1 Usage](#)
- [2 Functions](#)
- [3 Data modules](#)
- [4 Copyright](#)

Usage

This module provides functions that access information on Unicode code points. The information is retrieved from data modules generated from the [Unicode Character Database](#), or derived by rules given in the [Unicode Specification](#). It and its submodules were copied from English Wiktionary and then modified; see [there](#) for more information.

Functions

[Vorlage:Code](#)

Receives a code point (number) and returns its name or label; for example, [Vorlage:Code](#) returns [Vorlage:Code](#).

For example, [Vorlage:Tnull](#) <reserved-0061>

[Vorlage:Code](#)

Template-invokable functions that allow access to the functions starting with `lookup` and `is`. Replace the first underscore in the function name with a pipe. For most of the functions, add the code point in hexadecimal base as the next parameter, but for `is_Latin`, `is_rtl`, and `is_valid_pagename`, add text. [HTML character references](#) in the text are decoded by the module into code points.

For example, [Vorlage:Tnull](#) **Lua-Fehler in Zeile 293: attempt to index local 'data_module' (a boolean value).**

Data modules

The data used by functions in this module is found in [submodules](#). Some are generated by [AWK](#) scripts shown at [User:Kephir/Unicode](#) on English Wiktionary, others by Lua scripts on the `/make` subpages of the submodules.

- [Module:Unicode data/aliases](#): the formal name aliases for characters (from [NameAliases.txt](#))
- [Module:Unicode data/blocks](#): the list of Unicode blocks (from [Blocks.txt](#))
- [Module:Unicode data/category](#): data mapping characters to their General Category (from [DerivedGeneralCategory.txt](#))
- [Module:Unicode data/control](#): data for identifying characters that belong to the General Categories of Separator and Other (from [DerivedGeneralCategory.txt](#))
- [Module:Unicode data/combining](#): data mapping characters to their Combining Classes (from [DerivedCombiningClass.txt](#))
- [Module:Unicode data/Hangul](#): data used to generate the names of [Hangul](#) syllables (from [Jamo.txt](#))
- [Module:Unicode data/scripts](#): data mapping characters to their Unicode script properties (from [Scripts.txt](#)).

The name data modules ([Module:Unicode data/names/xxx](#)) were compiled from [UnicodeData.txt](#). Each one contains, at maximum, code points U+xxx000 to U+xxxFFF. **Lua-Fehler in mw.title.lua, Zeile 206: too many expensive function calls**

Copyright

The Unicode database is released by Unicode Inc. under the following terms:

Copyright © 1991-2018 Unicode, Inc. All rights reserved. Distributed under the Terms of Use in <https://www.unicode.org/copyright.html>.

Permission is hereby granted, free of charge, to any person obtaining a copy of the Unicode data files and any associated documentation (the "Data Files") or Unicode software and any associated documentation (the "Software") to deal in the Data Files or Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, and/or sell copies of the Data Files or Software, and to permit persons to whom the Data Files or Software are furnished to do so, provided that either (a) this copyright and permission notice appear with all copies of the Data Files or Software, or (b) this copyright and permission notice appear in associated Documentation.

THE DATA FILES AND SOFTWARE ARE PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT OF THIRD PARTY RIGHTS. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR HOLDERS INCLUDED IN THIS NOTICE BE LIABLE FOR ANY CLAIM, OR ANY SPECIAL INDIRECT OR CONSEQUENTIAL DAMAGES, OR ANY DAMAGES WHATSOEVER RESULTING FROM LOSS OF USE, DATA OR PROFITS, WHETHER IN AN ACTION OF CONTRACT, NEGLIGENCE OR OTHER TORTIOUS ACTION, ARISING OUT OF OR IN CONNECTION WITH THE USE OR PERFORMANCE OF THE DATA FILES OR SOFTWARE.

Except as contained in this notice, the name of a copyright holder shall not be used in advertising or otherwise to promote the sale, use or other dealings in these Data Files or Software without prior written authorization of the copyright holder.

```
local p = {}

local floor = math.floor

local function errorf(level, ...)
    if type(level) == "number" then
        return error(string.format(...), level + 1)
    else -- level is actually the format string.
        return error(string.format(level, ...), 2)
    end
end

local function binary_range_search(codepoint, ranges)
    local low, mid, high
    low, high = 1, ranges.length or require "Module:TableTools".length(ranges)
    while low <= high do
        mid = floor((low + high) / 2)
        local range = ranges[mid]
        if codepoint < range[1] then
            high = mid - 1
        elseif codepoint <= range[2] then
            return range, mid
        else
            low = mid + 1
        end
    end
    return nil, mid
end
p.binary_range_search = binary_range_search

--[
local function linear_range_search(codepoint, ranges)
    for i, range in ipairs(ranges) do
        if range[1] <= codepoint and codepoint <= range[2] then
            return range
        end
    end
end

--]]

-- Load a module by indexing "loader" with the name of the module minus the
-- "Module:Unicode data/" part. For instance, loader.blocks returns
-- [[Module:Unicode data/blocks]]. If a module cannot be loaded, false will be
-- returned.
local loader = setmetatable({}, {
    __index = function (self, key)
        local success, data = pcall(mw.loadData, "Module:Unicode data/" .. key)
        if not success then
            data = false
        end
        self[key] = data
        return data
    end
})

-- For the algorithm used to generate Hangul Syllable names,
-- see "Hangul Syllable Name Generation" in section 3.12 of the
-- Unicode Specification:
-- https://www.unicode.org/versions/Unicode11.0.0/ch03.pdf
local name_hooks = {
```

```

{      0x00,      0x1F, "<control-%04X>" }, -- C0 control characters
{      0x7F,      0x9F, "<control-%04X>" }, -- DEL and C1 control characters
{      0x3400,    0x4DBF, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph Extension A
{      0x4E00,    0x9FFF, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph
{      0xAC00,    0xD7A3, function (codepoint) -- Hangul Syllables
    local Hangul_data = loader.Hangul
    local syllable_index = codepoint - 0xAC00

    return ("HANGUL SYLLABLE %s%s%s"):format(
        Hangul_data.leads[floor(syllable_index / Hangul_data.final_count)]
        Hangul_data.vowels[floor((syllable_index % Hangul_data.final_count
            / Hangul_data.trail_count)],
        Hangul_data.trails[syllable_index % Hangul_data.trail_count]
    )
end },
-- High Surrogates, High Private Use Surrogates, Low Surrogates
{      0xD800,    0xDFFF, "<surrogate-%04X>" },
{      0xE000,    0xF8FF, "<private-use-%04X>" }, -- Private Use
-- CJK Compatibility Ideographs
{      0xF900,    0xFA6D, "CJK COMPATIBILITY IDEOGRAPH-%04X" },
{      0xFA70,    0xFAD9, "CJK COMPATIBILITY IDEOGRAPH-%04X" },
{      0x17000,   0x187F7, "TANGUT IDEOGRAPH-%04X" }, -- Tangut Ideograph
{      0x18800,   0x18AFF, function (codepoint)
    return ("TANGUT COMPONENT-%03d"):format(codepoint - 0x187FF)
end },
{      0x18D00,   0x18D08, "TANGUT IDEOGRAPH-%04X" }, -- Tangut Ideograph Supplement
{      0x1B170,   0x1B2FB, "NUSHU CHARACTER-%04X" }, -- Nushu
{      0x20000,   0x2A6DF, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph Extension B
{      0x2A700,   0x2B738, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph Extension C
{      0x2B740,   0x2B81D, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph Extension D
{      0x2B820,   0x2CEA1, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph Extension E
{      0x2CEB0,   0x2EBE0, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph Extension F
-- CJK Compatibility Ideographs Supplement (Supplementary Ideographic Plane)
{      0x2F800,   0x2FA1D, "CJK COMPATIBILITY IDEOGRAPH-%04X" },
{      0xE0100,   0xE01EF, function (codepoint) -- Variation Selectors Supplement
    return ("VARIATION SELECTOR-%d"):format(codepoint - 0xE0100 + 17)
end},
{      0x30000,   0x3134A, "CJK UNIFIED IDEOGRAPH-%04X" }, -- CJK Ideograph Extension G
{      0xF0000,   0xFFFFFD, "<private-use-%04X>" }, -- Plane 15 Private Use
{      0x100000,  0x10FFFFD, "<private-use-%04X>" } -- Plane 16 Private Use
}
name_hooks.length = #name_hooks

local name_range_cache

local function generate_name(data, codepoint)
    if type(data) == "string" then
        return data:format(codepoint)
    else
        return data(codepoint)
    end
end

end

--[
-- Checks that the code point is a number and in range.
-- Does not check whether code point is an integer.
-- Not used
local function check_codepoint(funcName, argIdx, val)
    require 'libraryUtil'.checkType(funcName, argIdx, val, 'number')
    if codepoint < 0 or 0x10FFFF < codepoint then
        errorf("Codepoint %04X out of range", codepoint)
    end
end

end
--]]

-- https://www.unicode.org/versions/Unicode11.0.0/ch04.pdf, section 4.8
function p.lookup_name(codepoint)

```

```

-- U+FDD0-U+FDEF and all code points ending in FFFE or FFFF are Unassigned
-- (Cn) and specifically noncharacters:
-- https://www.unicode.org/faq/private_use.html#nonchar4
if 0xFDD0 <= codepoint and (codepoint <= 0xFDEF
    or floor(codepoint % 0x10000) >= 0xFFFFE) then
    return ("

```

```

-- An ipairs-type iterator generator for the list of blocks.
function p.enum_blocks()
    local blocks = loader.blocks
    return block_iter, blocks, 0
end

function p.lookup_plane(codepoint)
    local i = floor(codepoint / 0x10000)
    return planes[i] or ("Plane %u"):format(i)
end

function p.lookup_block(codepoint)
    local blocks = loader.blocks
    local range = binary_range_search(codepoint, blocks)
    if range then
        return range[3]
    else
        return "No Block"
    end
end

function p.get_block_info(name)
    for i, block in ipairs(loader.blocks) do
        if block[3] == name then
            return block
        end
    end
end

function p.is_valid_pagename(pagename)
    local has_nonws = false

    for cp in mw.ustring.gcodepoint(pagename) do
        if (cp == 0x0023) -- #
        or (cp == 0x005B) -- [
        or (cp == 0x005D) -- ]
        or (cp == 0x007B) -- {
        or (cp == 0x007C) -- |
        or (cp == 0x007D) -- }
        or (cp == 0x180E) -- MONGOLIAN VOWEL SEPARATOR
        or ((cp >= 0x2000) and (cp <= 0x200A)) -- spaces in General Punctuation block
        or (cp == 0xFFFFD) -- REPLACEMENT CHARACTER
        then
            return false
        end

        local printable, result = p.is_printable(cp)
        if not printable then
            return false
        end

        if result ~= "space-separator" then
            has_nonws = true
        end
    end

    return has_nonws
end

local function manual_unpack(what, from)
    if what[from + 1] == nil then
        return what[from]
    end

    local result = {}
    from = from or 1
    for i, item in ipairs(what) do

```

```

        if i >= from then
            table.insert(result, item)
        end
    end
    return unpack(result)
end

local function compare_ranges(range1, range2)
    return range1[1] < range2[1]
end

-- Creates a function to look up data in a module that contains "singles" (a
-- code point-to-data map) and "ranges" (an array containing arrays that contain
-- the low and high code points of a range and the data associated with that
-- range).
-- "loader" loads and returns the "singles" and "ranges" tables.
-- "match_func" is passed the code point and either the data or the "dots", and
-- generates the final result of the function.
-- The varargs ("dots") describes the default data to be returned if there wasn't
-- a match.
-- In case the function is used more than once, "cache" saves ranges that have
-- already been found to match, or a range whose data is the default if there
-- was no match.
local function memo_lookup(data_module_subpage, match_func, ...)
    local dots = { ... }
    local cache = {}
    local singles, ranges

    return function (codepoint)
        if not singles then
            local data_module = loader[data_module_subpage]
            singles, ranges = data_module.singles, data_module.ranges
        end

        if singles[codepoint] then
            return match_func(codepoint, singles[codepoint])
        end

        local range = binary_range_search(codepoint, cache)
        if range then
            return match_func(codepoint, manual_unpack(range, 3))
        end

        local range, index = binary_range_search(codepoint, ranges)
        if range then
            table.insert(cache, range)
            table.sort(cache, compare_ranges)
            return match_func(codepoint, manual_unpack(range, 3))
        end

        if ranges[index] then
            local dots_range
            if codepoint > ranges[index][2] then
                dots_range = {
                    ranges[index][2] + 1,
                    ranges[index + 1] and ranges[index + 1][1] - 1 or 0,
                    unpack(dots)
                }
            else -- codepoint < range[index][1]
                dots_range = {
                    ranges[index - 1] and ranges[index - 1][2] + 1 or 0,
                    ranges[index][1] - 1,
                    unpack(dots)
                }
            end
            return match_func(codepoint, dots_range)
        end
        table.sort(cache, compare_ranges)
    end
end

```

```

        return match_func(codepoint)
    end
end

-- Get a code point's combining class value in [[Module:Unicode data/combining]],
-- and return whether this value is not zero. Zero is assigned as the default
-- if the combining class value is not found in this data module.
-- That is, return true if character is combining, or false if it is not.
-- See https://www.unicode.org/reports/tr44/#Canonical_Combining_Class_Values for
-- more information.
p.is_combining = memo_lookup(
    "combining",
    function (codepoint, combining_class)
        return combining_class and combining_class ~= 0 or false
    end,
    0)

function p.add_dotted_circle(str)
    return (mw.usttring.gsub(str, ".",
        function(char)
            if p.is_combining(mw.usttring.codepoint(char)) then
                return ' ' .. char
            end
        end))
end

local lookup_control = memo_lookup(
    "control",
    function (codepoint, ccc)
        return ccc or "assigned"
    end,
    "assigned")
p.lookup_control = lookup_control

function p.is_assigned(codepoint)
    return lookup_control(codepoint) ~= "unassigned"
end

function p.is_printable(codepoint)
    local result = lookup_control(codepoint)
    return (result == "assigned") or (result == "space-separator"), result
end

function p.is_whitespace(codepoint)
    local result = lookup_control(codepoint)
    return (result == "space-separator"), result
end

p.lookup_category = memo_lookup(
    "category",
    function (codepoint, category)
        return category
    end,
    "Cn")

local lookup_script = memo_lookup(
    "scripts",
    function (codepoint, script_code)
        return script_code or 'Zzzz'
    end,
    "Zzzz")
p.lookup_script = lookup_script

function p.get_best_script(str)
    -- Check type of argument, because mw.text.decode coerces numbers to strings!
    require "libraryUtil".checkType("get_best_script", 1, str, "string")

```

```

-- Convert HTML character references (including named character references,
-- or character entities) to characters.
str = mw.text.decode(str, true)

local scripts = {}
for codepoint in mw.ustr.gcodepoint(str) do
    local script = lookup_script(codepoint)

    -- Ignore "Inherited", "Undetermined", or "Uncoded" scripts.
    if not (script == "Zyyy" or script == "Zinh" or script == "Zzzz") then
        scripts[script] = true
    end
end

-- If scripts does not contain two or more keys,
-- return first and only key (script code) in table.
if not next(scripts, next(scripts)) then
    return next(scripts)
end -- else return majority script, or else "Zzzz"?

end

function p.is_Latin(str)
    require "libraryUtil".checkType("get_best_script", 1, str, "string")
    str = mw.text.decode(str, true)

    -- Search for the leading bytes that introduce the UTF-8 encoding of the
    -- code points U+0340-U+10FFFF. If they are not found and there is at least
    -- one Latin-script character, the string counts as Latin, because the rest
    -- of the characters can only be Zyyy, Zinh, and Zzzz.
    -- The only scripts found below U+0370 (the first code point of the Greek
    -- and Coptic block) are Latn, Zyyy, Zinh, and Zzzz.
    -- See the codepage in the [[UTF-8]] article.
    if not str:find "[\205-\244]" then
        for codepoint in mw.ustr.gcodepoint(str) do
            if lookup_script(codepoint) == "Latn" then
                return true
            end
        end
    end

    local Latn = false

    for codepoint in mw.ustr.gcodepoint(str) do
        local script = lookup_script(codepoint)

        if script == "Latn" then
            Latn = true
        elseif not (script == "Zyyy" or script == "Zinh"
            or script == "Zzzz") then
            return false
        end
    end

    return Latn
end

-- Checks that a string contains only characters belonging to right-to-left
-- scripts, or characters of ignorable scripts.
function p.is_rtl(str)
    require "libraryUtil".checkType("get_best_script", 1, str, "string")
    str = mw.text.decode(str, true)

    -- Search for the leading bytes that introduce the UTF-8 encoding of the
    -- code points U+0580-U+10FFFF. If they are not found, the string can only
    -- have characters from a left-to-right script, because the first code point
    -- in a right-to-left script is U+0591, in the Hebrew block.

```

```

    if not str:find "[\214-\244]" then
        return false
    end

    local result = false
    local rtl = loader.scripts.rtl
    for codepoint in mw.ustr.gcodepoint(str) do
        local script = lookup_script(codepoint)

        if rtl[script] then
            result = true
        elseif not (script == "Zyyy" or script == "Zinh"
            or script == "Zzzz") then
            return false
        end
    end

    return result
end

local function get_codepoint(args, arg)
    local codepoint_string = args[arg]
        or errorf(2, "Parameter %s is required", tostring(arg))
    local codepoint = tonumber(codepoint_string, 16)
        or errorf(2, "Parameter %s is not a code point in hexadecimal base",
            tostring(arg))
    if not (0 <= codepoint and codepoint <= 0x10FFFF) then
        errorf(2, "code point in parameter %s out of range", tostring(arg))
    end
    return codepoint
end

local function get_func(args, arg, prefix)
    local suffix = args[arg]
        or errorf(2, "Parameter %s is required", tostring(arg))
    suffix = mw.text.trim(suffix)
    local func_name = prefix .. suffix
    local func = p[func_name]
        or errorf(2, "There is no function '%s'", func_name)
    return func
end

-- This function allows any of the "lookup" functions to be invoked. The first
-- parameter is the word after "lookup_"; the second parameter is the code point
-- in hexadecimal base.
function p.lookup(frame)
    local func = get_func(frame.args, 1, "lookup_")
    local codepoint = get_codepoint(frame.args, 2)
    local result = func(codepoint)
    if func == p.lookup_name then
        -- Prevent code point labels such as <control-0000> from being
        -- interpreted as HTML tags.
        result = result:gsub("<", "&lt;")
    end
    return result
end

function p.is(frame)
    local func = get_func(frame.args, 1, "is_")

    -- is_Latin and is_valid_pagename take strings.
    if func == p.is_Latin or func == p.is_valid_pagename or func == p.is_rtl then
        return (func(frame.args[2]))
    else -- The rest take code points.
        local codepoint = get_codepoint(frame.args, 2)
        return (func(codepoint)) -- Adjust to one result.
    end
end

```

end

return p