

Proposed Draft Unicode® Standard Annex #60

DATA FOR NON HAN IDEOGRAPHIC SCRIPTS

Version	Unicode 18.0.0
Editor	Michel Suignard
Date	2026-02-04
This Version	https://www.unicode.org/reports/tr60/tr60-2.html
Previous Version	https://www.unicode.org/reports/tr60/tr60-1.html
Latest Version	https://www.unicode.org/reports/tr60/
Latest Proposed Update	https://www.unicode.org/reports/tr60/proposed.html
Revision	2

Summary

*This document describes the Sources and other ancillary data for non Han Ideographic Scripts, including Jurchen, Nüshu, **Seal**, and Tangut.*

Status

*This is a **draft** document which may be updated, replaced, or superseded by other documents at any time. Publication does not imply endorsement by the Unicode Consortium. This is not a stable document; it is inappropriate to cite this document as other than a work in progress.*

A Unicode Standard Annex (UAX) forms an integral part of the Unicode Standard, but is published online as a separate document. The Unicode Standard may require conformance to normative content in a Unicode Standard Annex, if so specified in the Conformance chapter of that version of the Unicode Standard. The version number of a UAX document corresponds to the version of the Unicode Standard of which it forms a part.

Please submit corrigenda and other comments with the online reporting form [[Feedback](#)]. Related information that is useful in understanding this annex is found in Unicode Standard Annex #41, “[Common References for Unicode Standard Annexes](#).” For the latest version of the Unicode Standard, see [[Unicode](#)]. For a list of current Unicode Technical Reports, see [[Reports](#)]. For more information about versions of the Unicode Standard, see [[Versions](#)]. For any errata which may apply to this annex, see [[Errata](#)].

Contents

- 1 Introduction
 - 2 Mechanics
 - 2.1 Data files Design
 - 2.2 Data files for Jurchen, Nüshu, Seal, and Tangut
 - 3 Property Types
 - 3.1 Sources
 - 3.2 Radical-Stroke Counts
 - 3.3 Readings
 - 3.4 Numeric values
 - 3.5 Other Data
 - 4 Scripts Properties
 - 4.1 Jurchen
 - 4.2 Nüshu
 - 4.3 Seal
 - 4.4 Tangut
 - 5 History
 - Acknowledgements
 - Modifications
-

1 Introduction

This document is a guide to information including sources and other ancillary data related to ideographic scripts other than Han. Historically, a summary and often incomplete version of that information was provided in the data file preambles related to these scripts. This document formalizes these elements in a structure similar to what is done for Han characters in UAX #38 Unihan Han Database information. In common with Han ideographs, elements of these other ideographic scripts are encoded using algorithmic names, including the name of the script and a multi-digit notation indicating the hexadecimal value of the code point, therefore providing little information about the identity of the character. The ancillary data provided by the related data files define additional information such as the various sources for the ideograph identity, and other ancillary information, such as the reading and radical-stroke index. While sources are always provided, the ancillary information varies between scripts. Similar to the Unihan database, this information could grow in the future, such as adding sources or other types of data related to specific code points.

The scripts covered by this document include Jurchen, Nüshu, Seal, and Tangut, referred as 'covered scripts' in the rest of this document. Note that while another East Asian encoded script, Khitan Small Script, had properties documented in the various encoding proposals, especially <https://www.unicode.org/L2/L2016/16113r-n4725r-khitan-small-script.pdf>, they were not surfaced in any Unicode data files. Nothing precludes their addition in the future, as it would improve the knowledge related to that script.

This document is a guide to these data files, one per covered script, describing their mechanics, the nature of their contents, and the status of the various properties. One the main goal of this document is to provide a single point of reference for all property information related to the covered scripts.

2 Mechanics

2.1 Data files Design

The data files consist of a number of fields containing data for each of the covered script's ideographs included in the Unicode Standard. The fields, all of which correspond to

properties, have names that consist entirely of ASCII letters and digits with no spaces or other punctuation except for underscore. For historical reasons, they all start with a lowercase k.

All data in these data files is stored in UTF-8 using Normalization Form C (NFC). Note, however, that the “Syntax” descriptions below, used for validation of property values, operate on Normalization Form D (NFD), primarily because that makes the regular expressions simpler.

2.2 Data files for Jurchen, Nüshu, Seal, and Tangut

Included with the [UCD] are three four files called `JurchenSources.txt`, `NushuSources.txt`, `SealSources.txt`, and `TangutSources.txt`. These files are single text files, in UTF-8, NFC, and using Unix line endings which contain the values for all properties related to each of the covered scripts. Properties are described by categories in this document but are nevertheless included in a single file per script (unlike, for example the Unihan database which is made of multiple files for the Han script). For most scripts, the All properties use a 'k' prefix followed by the four-letter abbreviated version of the script name as described in `PropertyValueAliases.txt`. For example, for the Jurchen script, the prefix is `kJURC`, and an example of property value is `kJURC_Src`. The notable exception is the Tangut script where the original 'TGT' abbreviation has been kept.

Review Note: There is ongoing work to clarify the status of these data sets in term of Unicode properties. These data files currently contain data which is only in scope for the script they are addressing. In that aspect they are different from typical Unicode properties which encompass the whole Unicode repertoire. As such, they were not subject to the typical constraints of Unicode properties, such stability, consistency, etc. By moving these data definitions in a UAX, the use of their status as Normative or Informative creates a stability requirement that may not be desired. At this moment, only properties that are essential to identities are qualified as 'normative', all others are qualified as 'provisional'. However, it may be desirable to make all properties of type 'Radical-Stroke-Counts' informative to be consistent with the similar Unihan property. This may also require an update in UAX #44 concerning the description of these properties.

In this file, blank lines may be ignored; lines beginning with # are comment lines used to provide the header and footer. Each of the remaining lines is one entry, with three tab-separated fields: the Unicode Scalar Value, the property name, and the value for the property for the given Unicode Scalar Value. For most of the properties, if multiple values are possible, the values are separated by spaces. No ideograph may have more than one instance of a given property associated with it, and no empty properties are included in these data files.

There is no formal limit on the lengths of any of the property values. Any Unicode character may be used in the property values except for control characters (especially tab, newline, and carriage return).

The data lines are sorted by Unicode Scalar Value and property-type as primary and secondary keys, respectively.

The file's header includes a summary of the properties each of these data files contains.

3 Property Types

The data in these data files serves a multitude of purposes, and the properties are grouped into categories according to the purpose they fulfill. A general discussion of the various categories is provided here, followed by a detailed description of the individual properties, alphabetically arranged. Among these categories, because the source information is essential in determining identity for characters which have algorithmically constructed names, the status of source related properties is 'normative'; all other properties have a 'provisional status'.

3.1 Sources

Sources are among the normative parts of these data files and refer to ideograph collections which identify encoded characters. These sources are defined as `kJURC_Src` for Jurchen, `kNSHU_DubenSrc` for Nüshu, `kSEAL_CCZSrc`, `kSEAL_DYCSrc`, `kSEAL_QJZSrc`, `kSEALTHXSrc` for Seal, and `kTGT_MergedSrc` for Tangut. These sources were typically documented in the encoding proposals for these scripts. Detailed descriptions of the syntax used for these sources are to be found in [Section 4, Script Properties](#), below.

3.2 Radical-Stroke Counts

Two of the scripts include radical-stroke counts: Jurchen with `kJURC_RSUnicode` and Tangut with `kTGT_RSUnicode`. All the radical-stroke properties used here are loosely derived from the radical system introduced by the 18th-century *Kangxi Dictionary* and used in the Unihan database for the Han ideographs. Each Tangut ideograph is assigned one of the 883 Tangut components, and each Jurchen ideographs is assigned one of the 51 Jurchen radicals. In all these cases, unlike Han, the component or radical assignment is never a semantic signifier; it is solely based on the ideograph's structure and is mainly meant to facilitate lookup of a specific ideograph in these large lists. The same two scripts also include a stroke count, and unlike the Han equivalent, the count includes the component or radical. It should be noted that the Nüshu repertoire is ordered by stroke counts (one to sixteen) but this is not reflected in any property. Finally, the Seal script specifies a Radical entry `kSEAL_Rad` which does not include a count and therefore is not included in this category.

3.3 Readings

Two of the scripts include a reading property: Jurchen with `kJURC_NCRreading` and Nüshu with `kNSHU_Reading`. Any attempt at providing a reading or set of readings for an ideograph is bound to be fraught with difficulty, because the readings will vary over time and from place to place, even within a language. However, because these readings have been documented in the encoding proposals and related to well-known sources, these are provided when available in these data files.

3.4 Numeric values

Two of the scripts include a numeric property: Jurchen with `kJURC_Numeric` and Tangut with `kTGT_Numeric`. This only applies to a few characters.

3.5 Other Data

This category includes properties that are typically specific to a given script. Currently only one property is two properties are defined in this category: `kJURC_Numeric`, `kSEAL_MCJK` and `kSEAL_Rad` used for Jurchen Seal.

4 Scripts Properties

Below is a listing of all properties for each of the covered scripts. Each of these lists is ordered alphabetically, with information on the property contents and syntax.

For each property we give the following information in the alphabetical listing: its *Property* tag, its *Unicode Status*, its *Category* as defined above, the Unicode version in which it was *Introduced*, its *Delimiter*, its *Syntax*, and its *Description*.

The *Property* name is the tag used in the data files to mark instances of this property.

The *Unicode Status* is either *Normative*, *Informative*, or *Provisional*, depending on whether it is a normative part of the standard, an informative part of the standard, or neither. We may also include *Deprecated* as a Unicode Status if the property is no longer to be used.

Properties which allow multiple property values have a *Delimiter* defined as “space” (U+0020 SPACE). Properties which do not have multiple property values have this defined as “N/A.” Some properties do not currently have multiple values in the data but may do so in the future.

For most properties with multiple values, the order of the values is arbitrary and has no particular significance. The most common order in such cases is alphabetical or numerical.

Validation is done as follows: The entry is split into subentries using the *Delimiter* (if defined), and each subentry converted to Normalization Form D (NFD). The value is valid if and only if each normalized subentry matches the property’s *Syntax* regular expression. Note that any given property’s *Syntax* is not guaranteed to be stable and may change in the future.

Finally, the *Description* contains not only a description of what the property contains, but also source information, known limitations, methodology used in deriving the data, and so on.

4.1 Jurchen

The properties covered in the table are: [kJURC_NCRReading](#), [kJURC_Numeric](#), [kJURC_RSUnicode](#), and [kJURC_Src](#).

Property	kJURC_NCRReading
Status	Provisional
Category	Readings
Introduced	18.0
Delimiter	N/A
Syntax	[^\t"]+
Default	N/A
Description	Reading given in Nǔzhēnwén Cídiǎn (Jīn), although it can be expressed in any Unicode character the value is typically a single string of Latin characters with optional parenthesis

Property	kJURC_Numeric
Status	Provisional
Category	Other Data: Numeric values
Introduced	18.0

Delimiter	N/A
Syntax	[1-9]\d{0,4}
Default	N
Description	Numeric value of the Jurchen character. It only applies to a few characters.

Property	kJURC_RSUnicode
Status	Provisional
Category	Sources ; Radical-Stroke
Introduced	18.0
Delimiter	N/A
Syntax	[1-9]\d{0,1}\.[1-9]\d{0,1}
Default	N/A
Description	The first number is the radical number, and the second number is the total stroke count.

Property	kJURC_Src
Status	Normative
Category	Description ; Sources
Introduced	18.0
Delimiter	N/A
Syntax	NC:\d{3}\.\d{2}(,\d{3}\.\d{2})? SJ-B:\d{3}[A-Z]\.\d JJ:\d{3} N5131\X-\d{4}
Default	N/A
Description	<p>The Jurchen sources are made of the following categories:</p> <p>NC Jīn Qǐzōng 金啓琮, Nǚzhēnwén Cídiǎn 女真文辞典 (Beijing: Wenwu chubanshe, 1984). The first number is the page number in Nǚzhēnwén Cídiǎn, the second number is the order of the entry on that page. There are multiple entries for some characters in the NC source, but this document only references a single entry for each character.</p> <p>SJ-B Berlin copy of the Sino-Jurchen Vocabulary. The first number is the folio, the second number is the position in # the folio.</p> <p>JJ Jin Guangping 金光平 and Jīn Qǐzōng 金啓琮, "Nǚzhen Yuyan Wenzī Yanjiū" 真语言文字研究 (Beijing: Wenwu chubanshe, 1980). The number indicates the page reference.</p> <p>N5131-X Sun Bojun, Nie Hongyin, Jing Yongshi, "A Supplementary Proposal to Encode the Jurchen Characters in UCS" (WG2 N5131), The sequence number is defined in WG2 N5131.</p>

4.2 Nūshu

The properties covered in the table are: [kNSHU_DubenSrc](#) and [kNSHU_Reading](#).

Property	kNSHU_DubenSrc
Status	Normative
Category	Sources
Introduced	10.0
Delimiter	N/A
Syntax	[1-9]\d.\d{2}
Default	N/A
Description	The only source documented in the file is Nǚshū Dúběn (NSDB) 女书读本 'Nüshu Reader': the first number is the page number in the NDSB, the second number is the order of the item on that page. While other sources have been mentioned in discussion about the proposal such as Nüshu Yongzi Bijiao[NSYZBJ] 女书用字比较 "A Comparison of characters used for writing Women's Script", they are not documented in the data file.

Property	kNSHU_Reading
Status	Provisional
Category	Readings
Introduced	10.0
Delimiter	N/A
Syntax	[a-z]+[1-9]\d{0,1}
Default	N/A
Description	Reading based on Nüshu Duben [NDSB], the numeric value after ascii text indicates the tones in five-degree contour tone marks

4.3 Seal

The properties covered in the table are: [kSEAL_CCZSrc](#), [kSEAL_DYCSrc](#), [kSEAL_MCJK](#), [kSEAL_QJZSrc](#), [kSEAL_Rad](#), and [kSEAL_THXSrc](#).

Property	kSEAL_CCZSrc
Status	Normative
Category	Sources
Introduced	18.0
Delimiter	N/A
Syntax	C-\d{5}
Default	N/A
Description	Chen Changzhi CCZ ((陳昌治單行本) source, one version of the Daxu Ben Shuowen Jiezi. The so-called "newly added characters (新附字)" specific to

	Daxu Ben versions are here numbered in the character sequence as 5-digit numbers.
--	---

Property	kSEAL_DYCSrc
Status	Normative
Category	Sources
Introduced	18.0
Delimiter	N/A
Syntax	D-\d{5}
Default	N/A
Description	Duan Zhu (DYC) source created by Duan Yucai (段玉裁).

Property	kSEAL_MCJK
Status	Provisional
Category	Other data
Introduced	18.0
Delimiter	N/A
Syntax	[0-9A-F]{4,5}
Default	N/A
Description	Each Seal character is associated with a single CJK Unified ideograph which may be used to refer to the Seal character. The value corresponds to the code point of the CJK Unified ideograph in hexadecimal notation. That ideograph value may be associated with multiple Seal characters.

Property	kSEAL_QJZSrc
Status	Normative
Category	Sources
Introduced	18.0
Delimiter	N/A
Syntax	K-\d{5}
Default	N/A
Description	Qi Junzao (QJZ) (祁雋藻刻本) source, one version of Xiaoxu Ben, work originally authored by Xu Kai (徐鍇)

Property	kSEAL_Rad
Status	Provisional
Category	Other data
Introduced	18.0
Delimiter	N/A
Syntax	\d{1,3}\.[0-9A-F]{5}
Default	N/A
Description	These values provide the radical number associated with each Seal character (1 to 540) along with the radical code point expressed in hexadecimal notation, the two values separated by a dot. In the Seal script, unlike many

	other scripts, the radicals are not encoded separately because they are just regular Seal characters. They are identified as being the first (and sometime unique) element of the group constituted of all elements sharing the same radical value.
--	---

Property	kSEAL_THXSrc
Status	Normative
Category	Sources
Introduced	18.0
Delimiter	N/A
Syntax	TH-(\d{5} X\d{3} Y\d{3})
Default	N/A
Description	Tenghuaxie THX (額勒 布藤花樹本) source, one version of the Daxu Ben Shuowen Jiezi. The so-called "newly added characters (新附字)" specific to Daxu Ben versions are here numbered in the character sequence as 3-digit numbers prefixed with 'X'. Because this source is also used to enumerate all encoded Seal characters, including characters not part of DaxuBen, those additional characters are referenced by 3-digit numbers prefixed by 'Y'.

4.4 Tangut

The properties covered in the table are: [kTGT_MergedSrc](#), [kTGT_Numeric](#), and [kTGT_RSUnicode](#).

<

Property	kTGT_MergedSrc
Status	Normative
Category	Description ; Sources
Introduced	9.0
Delimiter	N/A
Syntax	H2004-[AB]-\d{4} H2021-\d{6} L(19(86 97) 20(06 12))-\d{4} L2008-\d{4}([AB])-\d{4})? N1966-\d{3}-\d{2}[0-9A-Z]{1,2} N5217-\d{2} N5314-\d{2} S1968-\d{4} UTN42-\d{3}
Default	N/A
Description	The Tangut sources are made of the following categories: H2004 = Hán Xiǎománg (韓小忙), 西夏文正字研究 (Xīxiàwén Zhèngzì Yánjiū) [Research into the Correct Forms of Tangut Characters]. 2004. H2021 = Hán Xiǎománg (韓小忙), 西夏文词典: 世俗文献部分 (Xīxiàwén Cídiǎn: Shìsú Wénxiàn Bùfēn) [Tangut Word Dictionary: Secular Literature Part, 9 vols.]. 2021. WG2 N5286 2024-10-14.

	<p>L1986 = Lǐ Fànwén (李範文), 同音研究 (Tóngyīn Yánjiū) [Study of the Homophones]. Yinchuan. 1986.</p> <p>L1997 = Lǐ Fànwén (李範文), 夏漢字典 (Xià-Hàn Zìdiàn) [Tangut-Chinese Dictionary]. Beijing. 1997.</p> <p>L2006 = Lǐ Fànwén (李範文), 《五音切韻》与《文海宝韻》比较研究 (Wūyīn Qiēyùn yǔ Wénhǎi Bǎoyùn bǐjiào yánjiū), In 西夏研究 (Xīxià Yánjiū) [Western Xia Studies] no.2. Beijing. 2006</p> <p>L2008 = Lǐ Fànwén (李範文). 夏漢字典 (Xià-Hàn Zìdiàn) [Tangut-Chinese Dictionary]. Beijing, 2008.</p> <p>L2012 = Lǐ Fànwén, 2012 abridged edition, 2008 Tangut-Chinese Dictionary, cited in WG2 N 4724, page 2, 2014-04-21.</p> <p>N1966 = Nishida Tatsuo (西田龍雄), 西夏文小字典 (Seikabun Shōjiten) [Little Dictionary of Tangut], In 西夏語の研究 (Seikago no kenkyū) [A Study of the Hsi-Hsia Language] (1964-1966) vol.2. Tokyo, 1966.</p> <p>N5217 = Andrew West, Proposal to encode 2 Tangut components and 28 Tangut ideographs, WG2 N5217 = L2/23-149. 2023-10-02.</p> <p>N5314 = Andrew West, Proposal to encode one newly-identified Tangut ideograph, WG2 N5314 = L2/25-165. 2025-05-26.</p> <p>S1968 = Sofronov M. V. (М. В. Софронов), Грамматика тангутского языка (Grammatika tangutskogo jazyka) [Grammar of the Tangut Language]. Moscow, 1968.</p> <p>UTN42 = Andrew West and Viacheslav Zaytsev, Tangut Character Additions and Glyph Corrections, Unicode Technical Note #42. 2019-12-21.</p>
--	--

Property	kTGT_Numeric
Status	Provisional
Category	Numeric values
Introduced	18.0
Delimiter	N/A
Syntax	\d+(\.5)?
Default	N/A
Description	Numeric value of the Tangut character. It only applies to a few characters. The main bibliographic information for this property is: Jiǎ Chángyè 贾常业, ed. Xīxiàwén Zìdiǎn 西夏文字典 (Tangut Dictionary). Lanzhou: 甘肃文化出版社 (Gansu Culture Press), 2019.5, ISBN 978-7- 5490-1785-0, p. 995

Property	kTGT_RSUnicode
Status	Provisional
Category	Source Radical-Stroke
Introduced	9.0
Delimiter	N/A
Syntax	[1-9]\d{0,2}\.[1-9]\d{0,1}
Default	N/A

Description	The first number is the component number, and the second number is the total stroke count.
-------------	--

5 History

For two of the scripts, Nüshu, and Tangut, the information presented in this document used to be partially located in preambles attached to each of the related data files. These data files were created using data originally present in the encoding proposals and their updates. The data file for Tangut was incorporated in Unicode 9.0 and in Unicode 10.0 for Nüshu. It was augmented by details found in original encoding proposals for the covered scripts. Similarly, for the other two scripts, Jurchen and Seal, which are being encoded in Unicode 18.0, the information was directly extracted from the encoding proposals material.

References

For references for this annex, see Unicode Standard Annex #41, “[Common References for Unicode Standard Annexes](#).”

Acknowledgements

Andrew West (RIP) was the author of the encoding proposal for two scripts covered by this annex, Jurchen and Tangut, and he provided most of the information provided in this annex for these two scripts. For the other two scripts, Nüshu and Seal, the information gathering was a collective effort from many experts.

Modifications

Revision 2

- **Draft** of the first version of UAX#60 for Unicode 18.0.0.
- Changed the prefix for Tangut property names from 'TANG' to 'TGT'.
- Added the new kTGT_Numeric property.
- Added a new source for Tangut (N5314).
- Added coverage for the Seal script.

Revision 1

- **Proposed draft** of the first version of UAX#60 for Unicode 18.0.0.

Previous revisions will be accessed with the “Previous Version” link in the header when appropriate.

© 2024–2026 Unicode, Inc. This publication is protected by copyright, and permission must be obtained from Unicode, Inc. prior to any reproduction, modification, or other use not permitted by the [Terms of Use](#). Specifically, you may make copies of this publication and may annotate and translate it solely for personal or internal business purposes and not for public distribution, provided that any such permitted copies and modifications fully reproduce all copyright and other legal notices contained in the original. You may not make copies of or modifications to this publication for public distribution, or incorporate it in whole or in part into any product or publication without the express written permission of Unicode.

Use of all Unicode Products, including this publication, is governed by the Unicode [Terms of Use](#). The authors, contributors, and publishers have taken care in the preparation of this publication, but make no express or implied representation or warranty of any kind and assume no responsibility or liability for errors or omissions or for consequential or incidental damages that may arise therefrom. This publication is provided “AS-IS” without charge as a convenience to users.

Unicode and the Unicode Logo are registered trademarks of Unicode, Inc., in the United States and other countries.