

Note: This Research is created using Gemini 2.5 Pro.

Research Goals:

- (1) Find the latest official Unicode Standard documentation and chapter that specifically details the Khmer script (U+1780–U+17FF and U+19E0–U+19FF).
- (2) From the official documentation, identify and define the different classes of Khmer characters, such as base consonants, sub-consonants (using the 'coeng' character), dependent vowels, diacritics, and signs.
- (3) For each identified character class, compile a comprehensive list of the specific Unicode characters and their corresponding code point ranges.
- (4) Synthesize the canonical ordering rules for combining these character classes to form a valid Khmer grapheme cluster or syllable. This should be presented as a clear, sequential model.
- (5) Provide several detailed examples of Khmer words, breaking them down into their constituent Unicode characters and demonstrating the correct encoding sequence according to the official standard.

# The Canonical Encoding of Khmer Script: A Definitive Guide to the Unicode Standard

## Section 1: Introduction to the Khmer Script in the Unicode Standard

### 1.1. The Khmer Script: A Complex Abugida

The Khmer script, known in its native language as អក្សរខ្មែរ (*Āksâr Khmêr*), is the official writing system for the Khmer language of Cambodia. As a member of the Brahmic family of scripts, its lineage traces back to the Pallava script of Southern India, and it has been in continuous use since at least the early 7th century CE.<sup>1</sup> The script is an abugida, or alphasyllabary, a system in which consonant characters possess an inherent vowel sound. This inherent vowel is pronounced unless modified by a dependent vowel sign or silenced by a specific diacritic.<sup>1</sup> This fundamental structure, while common among South and Southeast Asian scripts, is implemented in Khmer with a unique complexity that poses significant challenges for digital representation.

A defining characteristic of the Khmer script is its two-series consonant system. Each of the 33 consonants used in the modern language belongs to one of two registers, historically corresponding to voiceless and voiced sounds but now distinguished by the vowel quality they impart.<sup>1</sup> These are commonly referred to as the a-series (e.g., ក *kâ*, with an inherent vowel of /ɑ:/) and the o-series (e.g., ខ *kô*, with an inherent vowel of /ɔ:/).<sup>1</sup> The series of the base consonant in a syllable dictates the pronunciation of any attached dependent vowel sign. For example, the vowel sign ា (U+17B6) produces an /a:/ sound when attached to an a-series consonant (e.g., កា *ka*) but an /iə/ sound when attached to an o-series consonant (e.g., ខា *kéa*).<sup>4</sup> This dual-register system effectively doubles the phonetic range of the available vowel signs, but it also introduces a layer of contextual dependency that a digital encoding model must faithfully represent.

Furthermore, Khmer orthography makes extensive use of consonant clusters, which are represented visually by stacking consonants. The first consonant of a cluster is written in its full form, while subsequent consonants (typically one, but sometimes two) are rendered in a

smaller, subscript form known as a ជើងអក្សរ (*cheung âksâr*, literally "foot of a letter").<sup>1</sup> This stacking mechanism is central to the script's visual and phonetic structure, allowing for the compact representation of complex syllable onsets. However, from an encoding perspective, this visual stacking requires a logical, linear sequence of characters, including a specific control character to trigger the subscript formation. The interplay between the base consonant, its series, any subscript consonants, and the dependent vowels creates a sophisticated orthographic system that demands a highly precise and unambiguous encoding standard.<sup>6</sup>

## 1.2. The Evolution of Khmer Encoding and the Principle of Canonical Equivalence

The Khmer script was first introduced into the Unicode Standard in Version 3.0, released in September 1999.<sup>7</sup> This initial encoding provided the necessary characters but left some aspects of their combination and ordering underspecified. This lack of a rigid, formal structure led to a significant and persistent problem in the digital ecosystem: encoding ambiguity. For many Khmer words and syllables, it became possible to create multiple different sequences of Unicode characters that would, in many rendering environments, produce the exact same visual output.<sup>9</sup> For example, a consonant shifter might be encoded before a subscript or after it, yet a font's substitution logic might render both sequences identically.

This ambiguity, while seemingly a minor technical issue, had severe practical consequences. It undermined the fundamental integrity of digital text, leading to user confusion, inconsistent data entry, and critical failures in text processing.<sup>9</sup> Searching for a word became unreliable, as a document might contain a visually identical string with a different underlying byte sequence, causing the search to fail. Sorting algorithms would incorrectly separate what appeared to be the same word. Spell checkers would flag one valid encoding as an error while accepting another. This situation placed an unreasonable burden on users, who were expected to possess a deep, technical understanding of Unicode's internal logic to type "correctly".<sup>9</sup>

In response to these systemic data integrity failures, the Unicode Consortium and technical experts in the field undertook a significant effort to establish a single, unambiguous encoding model. This work culminated in the development of **Unicode Technical Note #61 (UTN #61), "Khmer Encoding Structure,"** authored by Martin Hosken.<sup>11</sup> This document, now the authoritative guide, established a clear and prescriptive set of rules for character ordering. The entire model is founded on a core principle:

**"One Visual Form, One Encoding".**<sup>9</sup> This principle mandates that for any correctly formed Khmer orthographic syllable, there is one and only one valid sequence of Unicode characters to represent it.

This represents a deliberate and necessary philosophical shift in the standard's approach, moving from a model of representational flexibility to one of prescriptive rigidity. The goal is

no longer simply to provide the characters but to define a formal grammar that ensures data is consistent, verifiable, and machine-processable. This canonical model removes ambiguity, thereby guaranteeing that if two Khmer strings look the same, they are, in fact, the same at the byte level. This shift has profound implications for software development. The complexity of the canonical model means that the burden of ensuring correctness is transferred from the end-user to the software itself. Modern input method editors (IMEs), or keyboards, can no longer function as simple key-to-codepoint mappers. Instead, they must incorporate the full logic of the Khmer orthographic syllable model, acting as sophisticated engines that interpret a user's keystrokes and automatically normalize them into the single, valid, canonical sequence.<sup>11</sup> This ensures that all text created with a compliant system adheres to the standard, regardless of the user's technical expertise.

### 1.3. Authoritative Sources and Unicode Blocks

The definitive standard for character encoding is maintained by the Unicode Consortium, and all specifications are detailed in *The Unicode Standard*.<sup>13</sup> As of 2024, the latest major release is Version 16.0.<sup>15</sup> For the specific implementation and ordering rules of the Khmer script, the primary authoritative document is the aforementioned

**Unicode Technical Note #61 (UTN #61).**<sup>11</sup> This report is based on the principles and formal grammar defined within that technical note.

The characters required to write Khmer are located in two primary blocks within the Unicode Basic Multilingual Plane (BMP) <sup>1</sup>:

1. **Khmer:** This block, spanning the range U+1780 to U+17FF, contains the core character set, including all consonants, independent and dependent vowels, and various diacritics and signs used in the script.<sup>7</sup>
2. **Khmer Symbols:** This block, spanning U+19E0 to U+19FF, contains a set of characters used primarily for representing lunar dates in traditional Khmer calendars.<sup>1</sup>

Any software or system aiming for conformant processing of Khmer text must correctly implement the characters from these blocks according to the canonical ordering rules detailed in the subsequent sections of this report.

## Section 2: A Comprehensive Taxonomy of Khmer Characters

### 2.1. Character Classification Methodology

To understand the canonical ordering rules for Khmer, it is essential to first classify each

character not by its phonetic value or traditional grammatical category alone, but by its functional role within the Unicode orthographic syllable model. The model defined in UTN #61 deconstructs a syllable into a sequence of specific functional slots (e.g., Base, Coeng, Shifter, Vowel, Modifier). A character's classification determines which of these slots it can occupy, and therefore dictates its correct position in an encoded string. This functional taxonomy is the prerequisite for applying the ordering rules and is the foundation for building a valid Khmer grapheme cluster. The following table and descriptions provide a definitive classification for all relevant characters in the Khmer Unicode block.

## 2.2. The Master Character Taxonomy Table

The table below serves as a comprehensive reference, mapping every character in the U+1780–U+17FF block to its functional class within the canonical encoding model.

**Table 1: Comprehensive Khmer Character Taxonomy (U+1780–U+17FF)**

Functional Class	Sub-Class	Glyph	Codepoint	Unicode Name	Function and Notes
<b>Base</b>	Consonant	ក	U+1780	KHMER LETTER KA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ខ	U+1781	KHMER LETTER KHA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	គ	U+1782	KHMER LETTER KO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ឃ	U+1783	KHMER LETTER KHO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ង	U+1784	KHMER LETTER NGO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ច	U+1785	KHMER LETTER CA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ឆ	U+1786	KHMER LETTER CHA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ជ	U+1787	KHMER LETTER CO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ឈ	U+1788	KHMER	Forms the base

				LETTER CHO	of a syllable. O-series.
<b>Base</b>	Consonant	ញ	U+1789	KHMER LETTER NYO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ដ	U+178A	KHMER LETTER DA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ច	U+178B	KHMER LETTER TTHA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ឌ	U+178C	KHMER LETTER DO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ឍ	U+178D	KHMER LETTER TTHO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ណ	U+178E	KHMER LETTER NNO	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ត	U+178F	KHMER LETTER TA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ថ	U+1790	KHMER LETTER THA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ទ	U+1791	KHMER LETTER TO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ធ	U+1792	KHMER LETTER THO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ន	U+1793	KHMER LETTER NO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ប	U+1794	KHMER LETTER BA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ផ	U+1795	KHMER LETTER PHA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ព	U+1796	KHMER	Forms the base

				LETTER PO	of a syllable. O-series.
<b>Base</b>	Consonant	ភ	U+1797	KHMER LETTER PHO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ម	U+1798	KHMER LETTER MO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	យ	U+1799	KHMER LETTER YO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	រ	U+179A	KHMER LETTER RO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ល	U+179B	KHMER LETTER LO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	វ	U+179C	KHMER LETTER VO	Forms the base of a syllable. O-series.
<b>Base</b>	Consonant	ឃ	U+179D	KHMER LETTER SHA	Forms the base of a syllable. O-series. Rare, for Pali/Sanskrit.
<b>Base</b>	Consonant	ង	U+179E	KHMER LETTER SSO	Forms the base of a syllable. A-series. Rare, for Pali/Sanskrit.
<b>Base</b>	Consonant	ស	U+179F	KHMER LETTER SA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ហ	U+17A0	KHMER LETTER HA	Forms the base of a syllable. A-series.
<b>Base</b>	Consonant	ឡ	U+17A1	KHMER LETTER LA	Forms the base of a syllable. O-series. Has no subscript form.
<b>Base</b>	Consonant	អ	U+17A2	KHMER	Forms the base

				LETTER QA	of a syllable. A-series. Glottal stop; carrier for vowels.
<b>Base</b>	Independent Vowel	ត	U+17A5	KHMER INDEPENDENT VOWEL QI	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ត្ថ	U+17A6	KHMER INDEPENDENT VOWEL QII	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ឧ	U+17A7	KHMER INDEPENDENT VOWEL QU	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ឧ̌	U+17A8	KHMER INDEPENDENT Vowel QUK	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ឧ̍	U+17A9	KHMER INDEPENDENT VOWEL QUU	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ឧ̎	U+17AA	KHMER INDEPENDENT VOWEL QUUV	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	យ	U+17AB	KHMER INDEPENDENT VOWEL RY	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	យ̌	U+17AC	KHMER INDEPENDENT VOWEL RYY	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	យ̍	U+17AD	KHMER INDEPENDENT VOWEL LY	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	យ̎	U+17AE	KHMER INDEPENDENT VOWEL LYY	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ឯ	U+17AF	KHMER INDEPENDENT VOWEL QE	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ឯ̌	U+17B0	KHMER INDEPENDENT VOWEL QAI	Forms the base of a syllable.
<b>Base</b>	Independent	ឯ̍	U+17B1	KHMER	Forms the base



	Vowel			INDEPENDENT VOWEL QOO TYPE ONE	of a syllable.
<b>Base</b>	Independent Vowel	ᄒ	U+17B2	KHMER INDEPENDENT VOWEL QOO TYPE TWO	Forms the base of a syllable.
<b>Base</b>	Independent Vowel	ᄒ	U+17B3	KHMER INDEPENDENT VOWEL QAU	Forms the base of a syllable.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17B6	KHMER VOWEL SIGN AA	Modifies the inherent vowel of a base.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17B7	KHMER VOWEL SIGN I	Modifies the inherent vowel of a base. Renders above.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17B8	KHMER VOWEL SIGN II	Modifies the inherent vowel of a base. Renders above.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17B9	KHMER VOWEL SIGN Y	Modifies the inherent vowel of a base. Renders above.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17BA	KHMER VOWEL SIGN YY	Modifies the inherent vowel of a base. Renders above.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17BB	KHMER VOWEL SIGN U	Modifies the inherent vowel of a base. Renders below.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17BC	KHMER VOWEL SIGN UU	Modifies the inherent vowel of a base. Renders below.
<b>Vowel</b>	Dependent Vowel	ᄒ	U+17BD	KHMER VOWEL SIGN UA	Modifies the inherent vowel of a base. Renders below.
<b>Vowel</b>	Dependent	ᄒ	U+17BE	KHMER VOWEL	Modifies vowel.

	Vowel (Two-Part)			SIGN OE	Renders as two parts (left and above).
<b>Vowel</b>	Dependent Vowel (Two-Part)	្រ្ា	U+17BF	KHMER VOWEL SIGN YA	Modifies vowel. Renders as two parts (left and above).
<b>Vowel</b>	Dependent Vowel (Two-Part)	្រ្ិ	U+17C0	KHMER VOWEL SIGN IE	Modifies vowel. Renders as two parts (left and above).
<b>Vowel</b>	Dependent Vowel	្រ្ា	U+17C1	KHMER VOWEL SIGN E	Modifies the inherent vowel of a base. Renders to the left.
<b>Vowel</b>	Dependent Vowel	្រ្ា	U+17C2	KHMER VOWEL SIGN AE	Modifies the inherent vowel of a base. Renders to the left.
<b>Vowel</b>	Dependent Vowel	្រ្ា	U+17C3	KHMER VOWEL SIGN AI	Modifies the inherent vowel of a base. Renders to the left.
<b>Vowel</b>	Dependent Vowel (Two-Part)	្រ្ា	U+17C4	KHMER VOWEL SIGN OO	Modifies vowel. Renders as two parts (left and right).
<b>Vowel</b>	Dependent Vowel (Two-Part)	្រ្ា	U+17C5	KHMER VOWEL SIGN AU	Modifies vowel. Renders as two parts (left and right).
<b>Modifier</b>	Diacritic	្រ្ា	U+17C6	KHMER SIGN NIKAHIT	Combining sign for final nasalization (anusvara).
<b>Shifter</b>	Consonant Shifter	្រ្ា	U+17C9	KHMER SIGN MUUSIKATOAN	Shifts a consonant from o-series to a-series.

<b>Shifter</b>	Consonant Shifter	◌̃	U+17CA	KHMER SIGN TRIISAP	Shifts a consonant from a-series to o-series.
<b>Modifier</b>	Diacritic	◌̣	U+17CB	KHMER SIGN BANTOC	Shortens the preceding vowel sound.
<b>Robat</b>	Diacritic	◌̤	U+17CC	KHMER SIGN ROBAT	Represents a final 'r' sound. Has special ordering rules.
<b>Modifier</b>	Diacritic	◌̥	U+17CD	KHMER SIGN TOANDAKHIAT	Silences the consonant it is placed over.
<b>Modifier</b>	Diacritic	◌̦	U+17CE	KHMER SIGN KAKABAT	Rare sign used for emphasis.
<b>Modifier</b>	Diacritic	◌̧	U+17CF	KHMER SIGN AHSDA	Rare sign used for emphasis.
<b>Modifier</b>	Diacritic	◌̨	U+17D0	KHMER SIGN SAMYOK SANNYA	Modifies vowel pronunciation.
<b>Modifier</b>	Diacritic	◌̩	U+17D1	KHMER SIGN VIRIAM	Used in Sanskrit/Pali transliteration to silence inherent vowel.
<b>Coeng-Former</b>	Control Character	◌̪	U+17D2	KHMER SIGN COENG	Invisible character; renders the following Base as a subscript.
<b>Final</b>	Spacing Sign	◌̫	U+17C7	KHMER SIGN REAHMUK	Spacing sign for final aspiration (visarga).
<b>Final</b>	Spacing Sign	◌̬	U+17C8	KHMER SIGN YUUKALEAPINTU	Spacing sign indicating a short inherent vowel with glottal stop.
<b>Other</b>	Punctuation	។	U+17D4	KHMER SIGN KHAN	Period/full stop.

Other	Punctuation	្ក	U+17D5	KHMER SIGN BARIYOOSAN	End of a chapter or text.
Other	Punctuation	្ខ	U+17D6	KHMER SIGN CAMNUC PII KUUH	Colon.
Other	Repetition Sign	្គ	U+17D7	KHMER SIGN LEK TOO	Repeats the preceding word.
Other	Currency Symbol	៛	U+17DB	KHMER CURRENCY SYMBOL RIEL	The symbol for the Cambodian riel.
Other	Deprecated	្ឃ	U+17A3	KHMER INDEPENDENT VOWEL QAA	Deprecated since Unicode 4.0. Use of U+17A2 is preferred.
Other	Deprecated	្ង	U+17A4	KHMER INDEPENDENT VOWEL QAA	Deprecated since Unicode 5.2. Use of U+17A2 U+17B6 is preferred.
Other	Deprecated	្ច	U+17D3	KHMER SIGN BATHAMASAT	Originally for lunar dates. Use of U+19E0 is preferred.

## 2.3. Detailed Class Descriptions

- Base Characters:** These are the foundational elements of any Khmer orthographic syllable. A syllable must begin with exactly one character from this class.<sup>11</sup> This class includes the 33 standard consonants (U+1780–U+17A2) and the 17 independent vowels (U+17A5–U+17B3). While linguistically distinct, consonants and independent vowels behave identically as the syllable's structural anchor in the Unicode model. The consonant អ (U+17A2) is particularly important as it functions as a null consonant or glottal stop, often serving as a "carrier" for dependent vowels when a syllable begins with a vowel sound.<sup>3</sup>
- Subscript Formation (Coeng-Former):** The character KHMER SIGN COENG (U+17D2) is the single most important control character for Khmer encoding. It is invisible and has no glyph of its own. Its sole function is to signal to a rendering engine that the

immediately following Base character should be transformed into its subscript (coeng) form and attached to the preceding base consonant.<sup>6</sup> Every subscript in Khmer text is logically represented by this two-character sequence:

U+17D2 + Base Character.

- **Dependent Vowels (U+17B6–U+17C5):** These characters cannot stand alone and must be attached to a Base character (or a consonant cluster) to modify its inherent vowel sound. They are combining marks, but their visual rendering is complex. Some render above the base (e.g., ៊ U+17B7), some below (e.g., ្ក U+17BB), and some to the left (e.g., ្គ U+17C1). A crucial sub-category is the **two-part vowel** (or circumgraph), such as ្ដ (U+17BE), ្ឋ (U+17BF), ្ឌ (U+17C0), ្ឍ (U+17C4), and ្ណ (U+17C5).<sup>20</sup> Although these vowels render with glyph components on multiple sides of the base consonant, they are encoded as a single, atomic codepoint that must be placed *after* the entire consonant cluster in the logical string. The rendering engine is solely responsible for splitting the character into its constituent visual parts and positioning them correctly.
- **Consonant Shifters (U+17C9, U+17CA):** The KHMER SIGN MUUSIKATOAN (្ណ) and KHMER SIGN TRIISAP (្ត) are diacritics that alter the inherent series of a consonant or an entire consonant cluster.<sup>20</sup> MUUSIKATOAN shifts an o-series consonant to the a-series (e.g., ង *ngô* becomes ង្ណ *ngâ*). Conversely, TRIISAP shifts an a-series consonant to the o-series (e.g., ស *sâ* becomes ស្ត *sô*). Their placement in the encoded string is rigidly defined and is one of the most common sources of encoding errors in non-compliant systems.
- **Diacritics and Modifying Signs:** This is a broad category of combining marks that provide further phonetic modification.
  - The KHMER SIGN ROBAT (្ណ, U+17CC) is a unique diacritic representing a final /r/ sound from a preceding syllable. It has special ordering rules, placing it immediately after the base consonant, before any subscripts.<sup>4</sup>
  - Other important modifiers include KHMER SIGN NIKAHIT (្ណ, U+17C6) for final nasalization, KHMER SIGN BANTOC (្ណ, U+17CB) for vowel shortening, and KHMER SIGN TOANDAKHIAT (្ណ, U+17CD), which indicates that the consonant it is attached to is silent.<sup>4</sup>
- **Finals and Punctuation:** Unlike modifiers, which are non-spacing combining marks, Finals are spacing characters that appear at the end of a syllable. This class includes KHMER SIGN REAHMUK (្ណ, U+17C7) and KHMER SIGN YUUKALEAPINTU (្ណ, U+17C8).<sup>22</sup> Punctuation marks, such as the period-like KHMER SIGN KHAN (្ណ, U+17D4) and the section-ending KHMER SIGN BARIYOOSAN (្ណ, U+17D5), function similarly to their Western counterparts and are not part of the core syllable structure.<sup>19</sup>
- **Other Characters:** This category includes Khmer digits (U+17E0–U+17E9), the currency symbol (U+17DB), and several deprecated characters. U+17A3 and U+17A4 are deprecated independent vowels whose use is strongly discouraged in modern text in favor of combinations using the null consonant U+17A2.<sup>7</sup> Their inclusion in the standard

is for backward compatibility only.

## Section 3: The Structure of the Khmer Orthographic Syllable

### 3.1. The Syllable as the Atomic Unit of Encoding

The canonical model for Khmer Unicode is built not around individual characters but around the **orthographic syllable**. This unit is the smallest valid, self-contained combination of characters that forms a "grapheme cluster" in the context of the script.<sup>6</sup> A string of Khmer text is fundamentally a sequence of these syllables. The validity of a Khmer string is therefore determined by whether it can be parsed into a sequence of well-formed syllables according to a strict, formal grammar.

This grammatical approach is a cornerstone of the modern standard. It moves beyond simple lists of preferred orderings and provides a machine-testable model for correctness. Any sequence of Khmer codepoints can be validated against this grammar. If it conforms, it is canonical; if it deviates, it is malformed. This provides a powerful and unambiguous basis for software implementation. Developers can create parsers or state machines that definitively validate, normalize, and correct Khmer text, ensuring data integrity across all applications. This formal structure is what enables the principle of "One Visual Form, One Encoding" to be practically realized. The rules of this grammar also reveal the critical role of invisible control characters. Characters like KHMER SIGN COENG (U+17D2), and in more complex historical text, the Zero-Width Joiner (ZWJ, U+200D) and Zero-Width Non-Joiner (ZWNJ, U+200C), are not optional stylistic elements; they are essential grammatical operators that determine the structure and validity of a syllable.<sup>11</sup> Failure to handle these control characters correctly will result in non-canonical and often incorrectly rendered text.

### 3.2. The Canonical Syllable Structure for Modern Khmer

For the vast majority of modern Khmer text, including the national language of Cambodia and related minority languages, a simplified version of the full orthographic syllable model is sufficient. This structure, defined in UTN #61, provides a clear and robust template for encoding.<sup>11</sup> It can be expressed using a formal grammar notation:

Syllable = Base Robot? Coengs? Shifter? Vowel? Modifiers? Final?

Each component in this structure represents a functional slot that can be filled by characters of the corresponding class, as defined in Section 2. The ? indicates that the component is optional.

A detailed breakdown of the components is as follows:

- **Base:** This is the only mandatory component of a syllable. It must be a single character from the Base class, which includes all consonants (U+1780–U+17A2) and independent vowels (U+17A5–U+17B3).
- **Robat?:** An optional KHMER SIGN ROBAT (U+17CC). If present, it must immediately follow the Base.
- **Coengs?:** An optional sequence of one or two subscript consonants. Each subscript is encoded as a KHMER SIGN COENG (U+17D2) followed by a Base character.
- **Shifter?:** An optional consonant shifter, either KHMER SIGN MUUSIKATOAN (U+17C9) or KHMER SIGN TRIISAP (U+17CA). It must follow the entire consonant group (Base + Coengs).
- **Vowel?:** An optional dependent vowel (U+17B6–U+17C5). This includes single-codepoint representations of two-part vowels.
- **Modifiers?:** An optional sequence of up to two non-spacing modifying signs, such as NIKAHIT (U+17C6) or BANTOC (U+17CB).
- **Final?:** An optional spacing final sign, either REAHMUK (U+17C7) or YUUKALEAPINTU (U+17C8).

This structure provides a clear, linear path for encoding any modern Khmer syllable. For example, to encode a syllable with a base, a subscript, a vowel, and a modifier, the characters must be stored in precisely that order.

### 3.3. The Full Authoritative Structure for Historical and Complex Text

To provide comprehensive support for all forms of the Khmer script, including historical orthographies like Middle Khmer and complex Sanskrit/Pali loanwords, the Unicode standard defines a more elaborate and powerful syllable structure.<sup>11</sup> This full model is essential for developers of foundational technologies like text rendering engines and archives dealing with historical manuscripts. While more complex, it is built upon the same principles as the modern structure.<sup>11</sup>

The full authoritative structure is:

Syllable = Base Robat? Coengs? SubSyllable (FinalCoeng SubSyllable?)\*

The initial components (Base, Robat, Coengs) are the same as in the modern model. The key differences lie in the subsequent parts:

- **SubSyllable:** This is a sub-grouping defined as Shifter? Vowels? Modifiers? Final?. This structure is more flexible than the modern model.
- **Vowels?:** Note the plural. The full model allows for sequences of multiple dependent vowel characters, a feature found in Middle Khmer but not in the modern language. The grammar in UTN #61 specifies the valid combinations of these vowel signs.
- **(FinalCoeng SubSyllable?)\*:** This component, which can be repeated, introduces the concept of a **final coeng**. A FinalCoeng is a subscript consonant that represents the *end* of a phonetic syllable, a practice common in Middle Khmer. To distinguish this from

a standard initial coeng, it is encoded with a ZWJ (Zero-Width Joiner, U+200D) before the KHMER SIGN COENG (U.g., ZWJ + U+17D2 + Base). This explicit use of a control character prevents ambiguity and ensures that only intentionally formed final coengs are considered valid.

This full model demonstrates the depth of the Unicode standard's design, which aims to be both practical for modern usage and robust enough to support the script's entire historical and literary heritage. For most developers, the modern structure is sufficient, but an awareness of the full model is crucial for building truly comprehensive and future-proof text processing systems.

## **Section 4: Canonical Character Ordering: The Definitive Rules**

The following set of rules translates the formal syllable grammar into a prescriptive, step-by-step guide for implementation. Adherence to this sequence is mandatory for generating canonical Khmer Unicode text. Each rule defines the required position of a character class within the linear byte stream of an orthographic syllable.

### **4.1. Rule 1: The Base Character is Always First**

Every Khmer orthographic syllable must begin with a single character from the Base class. This character can be any of the 33 consonants (U+1780–U+17A2) or 17 independent vowels (U+17A5–U+17B3).<sup>11</sup> This character serves as the anchor to which all other components of the syllable are logically attached. No other character type, including vowel signs or diacritics, can precede the Base character within a syllable.

### **4.2. Rule 2: Robat (U+17CC) Follows the Base Immediately**

If the KHMER SIGN ROBAT (្ក, U+17CC) is present in a syllable, it must be encoded immediately after the Base character. It must precede any Coeng sequences, Shifter, Vowel, or other modifying signs.<sup>11</sup> Even though it is a diacritic, its unique phonetic role and historical usage grant it this privileged position in the encoding order.

### **4.3. Rule 3: Subscript Consonants (Coengs) are Encoded Sequentially**

Subscript consonants, or coengs, are encoded after the Base (and after Robat, if present).



Each subscript is represented by a two-character sequence: the invisible control character KHMER SIGN COENG (្ក, U+17D2) followed by the Base character that is to be subscripted. For example, the cluster ក្រ is encoded as:

U+1780 (ក) + U+17D2 (្ក) + U+178A (រ)

In Modern Khmer, a syllable can have up to two initial coengs (forming a three-consonant cluster).<sup>6</sup> In such cases, the

Coeng sequences are encoded in the order they appear visually, from top to bottom. For example, a cluster with a base, a first subscript, and a second subscript would be encoded: Base + Coeng1\_Sequence + Coeng2\_Sequence.

There is one critical exception to this visual ordering rule:

- **Special Sub-Rule for Subscript Ro (រ):** In a consonant cluster with two subscripts, if one of them is KHMER LETTER RO (រ, U+179A), its Coeng sequence (U+17D2 U+179A) must *always* be encoded second, as the final part of the consonant cluster, regardless of its visual or phonetic position.<sup>11</sup> This is a non-intuitive but mandatory rule for canonical encoding. For instance, in a cluster visually rendered as Base-Ro-Ka, the encoding must be Base + Coeng\_Ka\_Sequence + Coeng\_Ro\_Sequence.

#### 4.4. Rule 4: Consonant Shifters Follow the Complete Consonant Cluster

The consonant shifters, KHMER SIGN MUUSIKATOAN (្គ, U+17C9) and KHMER SIGN TRIISAP (្ឃ, U+17CA), must be encoded *after* the entire consonant cluster. This means they follow the Base character and all of its associated Coeng sequences.<sup>11</sup> Placing the shifter anywhere else (e.g., between the base and a coeng) is a common error in older systems and produces non-canonical text. The shifter applies its series-changing effect to the entire consonant cluster as a single unit.

This rule is complicated by a context-sensitive behavior known as "**downshifting**." Under certain conditions, a shifter's visual form and function can change. For example, a TRIISAP (្ឃ) placed over certain consonant clusters may render as a BANTOC (្គ) and a NIKAHIT (្ង). The precise rules for when downshifting occurs are complex and are defined in UTN #61 based on the series of the consonants in the cluster.<sup>11</sup>

To handle cases where this automatic behavior is not desired, the standard provides an explicit control mechanism:

- **Using ZWNJ (U+200C) to Control Downshifting:** A Zero-Width Non-Joiner (U+200C) can be placed immediately after a shifter to prevent it from downshifting. The use of ZWNJ is only canonical in contexts where downshifting would otherwise occur according to the rules in UTN #61. This provides an unambiguous way to encode visually similar but functionally distinct forms.<sup>11</sup>

## 4.5. Rule 5: Dependent Vowels Follow the Shifter

All dependent vowel signs (U+17B6–U+17C5) are encoded after the complete consonant group, which includes the Base, any Robat, all Coengs, and any Shifter. This is true even for vowels that render to the left of the base consonant, such as ្រ (U+17C1). The logical encoding order follows phonetic order, not visual order.

For two-part vowels like ្រ (U+17BE), which render with components on multiple sides of the consonant cluster, only the single codepoint for the vowel is encoded in this position.<sup>20</sup> The text shaping engine is responsible for correctly parsing this single codepoint and rendering its multiple glyphs in their appropriate visual positions around the preceding consonant cluster.

## 4.6. Rule 6: Modifying Signs and Finals are Encoded Last

Any remaining diacritics, classified as Modifiers, are encoded after the dependent vowel. These are non-spacing combining marks such as NIKAHIT (្រ, U+17C6) or TOANDAKHIAT (្រ, U+17CD). The canonical model permits up to two Modifiers in a sequence.<sup>11</sup> Finally, any spacing signs classified as Finals, such as REAHMUK (្រ, U+17C7) or YUUKALEAPINTU (្រ, U+17C8), are encoded at the very end of the syllable's character sequence.<sup>11</sup> This completes the construction of a valid, canonical Khmer orthographic syllable.

# Section 5: Practical Implementation and Illustrative Examples

This section provides a series of practical examples to demonstrate the application of the canonical ordering rules. Each example deconstructs a common Khmer word into its constituent orthographic syllables and provides a detailed breakdown of the Unicode character sequence required for its correct encoding.

## 5.1. Example 1: Consonant Cluster (្រ - Khmer)

The word ្រ (*Khmer*) demonstrates the fundamental mechanism for encoding a consonant cluster using the KHMER SIGN COENG. The word consists of a single syllable.

- **Rendered:** ្រ
- **Canonical Sequence:** U+1781 U+17D2 U+1798 U+17C2 U+179A
- **Breakdown:**

Codepoint	Glyph	Unicode Name	Role	Rule Applied
U+1781	្រ	KHMER LETTER	Base	Rule 1

		KHA		
U+17D2	្ក	KHMER SIGN COENG	Coeng-Former	Rule 3
U+1798	ម	KHMER LETTER MO	Coeng (forms ្ក)	Rule 3
U+17C2	្ថ	KHMER VOWEL SIGN AE	Vowel	Rule 5
U+179A	រ	KHMER LETTER RO	Final Consonant	N/A (Treated as a separate syllable in some models, but here as a final consonant)

Note that the vowel sign ្ថ (U+17C2), which renders to the left of the cluster, is encoded *after* the entire cluster (្ក + ្ក + ម), following Rule 5.

### 5.2. Example 2: Consonant Shifter (ផ្អ - Flute)

The word ផ្អ (pəy) illustrates the correct placement of a consonant shifter. The shifter modifies the base consonant ប (a-series *bâ*) to sound like *pâ*.

- **Rendered:** ផ្អ
- **Canonical Sequence:** U+1794 U+17C9 U+17B8
- **Breakdown:**

Codepoint	Glyph	Unicode Name	Role	Rule Applied
U+1794	ប	KHMER LETTER BA	Base	Rule 1
U+17C9	្គ	KHMER SIGN MUUSIKATOAN	Shifter	Rule 4
U+17B8	ា	KHMER VOWEL SIGN II	Vowel	Rule 5

The MUUSIKATOAN (U+17C9) is placed directly after the Base and before the Vowel, adhering strictly to Rule 4.

### 5.3. Example 3: Complex Cluster with Ro (ស្រុក - District/Village)

The word ស្រុក (*srōk*) is a single syllable containing a consonant cluster. It demonstrates the standard encoding of a subscript.

- **Rendered:** ស្រុក

- **Canonical Sequence:** U+179F U+17D2 U+179A U+17BB U+1780
- **Breakdown:**

Codepoint	Glyph	Unicode Name	Role	Rule Applied
U+179F	ស	KHMER LETTER SA	Base	Rule 1
U+17D2	្ក	KHMER SIGN COENG	Coeng-Former	Rule 3
U+179A	រ	KHMER LETTER RO	Coeng (forms ្ក)	Rule 3
U+17BB	ុ	KHMER VOWEL SIGN U	Vowel	Rule 5
U+1780	ក	KHMER LETTER KA	Final Consonant	N/A

## 5.4. Additional Examples for Advanced Cases

- **Two-Part Vowel (កើត - to be born):** The word កើត (*kaət*) uses the two-part vowel ើ (U+17BE).
  - **Rendered:** កើត
  - **Canonical Sequence:** U+1780 U+17BE U+178F
  - **Breakdown:**

Codepoint	Glyph	Unicode Name	Role	Rule Applied
U+1780	ក	KHMER LETTER KA	Base	Rule 1
U+17BE	ើ	KHMER VOWEL SIGN OE	Vowel	Rule 5
U+178F	ត	KHMER LETTER TA	Final Consonant	N/A

Despite rendering with parts on the left (ុ) and above (ើ), the vowel is represented by the single codepoint U+17BE, placed after the base consonant as required by Rule 5.

- **Use of Robat (ធម៌ - dharma):** The word ធម៌ (*thôrm*) demonstrates the correct placement of ROBAT.
  - **Rendered:** ធម៌
  - **Canonical Sequence:** U+1792 U+17CC U+1798
  - **Breakdown:**

Codepoint	Glyph	Unicode Name	Role	Rule Applied
U+1792	ធម៌	KHMER LETTER THO	Base	Rule 1
U+17CC	្ក	KHMER SIGN	Robat	Rule 2

		ROBAT		
U+1798	ឃ	KHMER LETTER MO	Base of second syllable	Rule 1

The ROBAT (U+17CC) is encoded immediately after the base consonant រ (U+1792), before the following consonant ឃ (U+1798), as mandated by Rule 2.

These examples illustrate that the logical, encoded order of Khmer characters follows a consistent and predictable phonetic-based sequence, which often differs from the final visual arrangement. Correct implementation requires strict adherence to this logical order.

## Section 6: Conclusion: Implications for the Digital Ecosystem

### 6.1. Summary of Canonical Principles

The modern Unicode standard for the Khmer script represents a decisive move towards a robust, unambiguous, and computationally sound model for digital text. This model is governed by two foundational principles. First is the mandate of **"One Visual Form, One Encoding,"** which eradicates the ambiguity of earlier specifications and ensures that any given correctly rendered Khmer string corresponds to one and only one underlying sequence of codepoints. This principle is the bedrock of data integrity for the Khmer language in the digital sphere.

Second is the definition of a formal, hierarchical **orthographic syllable structure**. Khmer text is not a simple linear sequence of characters but a structured composition of these syllables. The validity and canonical form of any text are determined by its conformance to a strict grammatical model. This model, detailed in Unicode Technical Note #61, provides a clear, machine-testable framework that dictates the precise ordering of all character types, from base consonants and vowels to invisible control characters and diacritics. Together, these principles provide the stability and predictability necessary for all forms of digital text processing.

### 6.2. Recommendations for Software Developers

The successful implementation of the Khmer script across the digital ecosystem depends on developers correctly applying these canonical principles. The following recommendations are directed at key areas of software development:

- **Rendering Engines:** Developers of text shaping engines (such as HarfBuzz or

platform-native APIs) must implement the full shaping logic described in UTN #61. This includes not only the basic substitution of subscripts and positioning of vowels but also the complex, context-sensitive rules for consonant shifter downshifting and the correct parsing of multipart vowels from single codepoints. Furthermore, rendering engines should provide clear visual feedback for malformed or non-canonical sequences, typically by inserting a dotted circle (U+25CC) to indicate an invalid character combination, thereby guiding users and developers toward creating correct text.<sup>10</sup>

- **Input Methods (IMEs):** The complexity of the canonical ordering rules means that the burden of correctness must shift from the user to the software. Developers of keyboard applications and other input methods for Khmer must implement normalization routines. These routines should intercept user keystrokes and, regardless of the physical typing order, intelligently assemble and output the single, canonical character sequence for the intended syllable. An IME that allows a user to generate non-canonical text is fundamentally broken.
- **Databases and Search Engines:** For searching and collation (sorting) to function reliably, all Khmer text must be normalized to its canonical form before being stored or indexed.<sup>11</sup> Database systems should include a canonical normalization step as part of their data ingestion and sanitization pipeline. Similarly, search engine indexers must normalize text to ensure that queries match documents regardless of any potential legacy encoding variations. Failure to do so will result in incomplete search results and incorrect alphabetical sorting, undermining the system's utility.<sup>26</sup>

### 6.3. Guidance for Content Creators and Digital Archivists

The long-term preservation and interoperability of Khmer digital heritage depend on adherence to the canonical standard. Content creators, from journalists to publishers, should ensure they are using modern, Unicode-compliant software and fonts that correctly implement the ordering rules.

For digital archivists and linguists working with existing corpora, it is highly recommended to develop or use validation tools to audit text for non-canonical sequences. Many older digital documents may contain text created with non-compliant systems. A one-time normalization process to convert these legacy encodings to the canonical form is a critical step in ensuring that these valuable linguistic resources remain accessible, searchable, and usable with future technologies.

### 6.4. The Future of Khmer on the Web

The establishment of a rigorous, unambiguous encoding standard does more than solve technical problems; it solidifies the foundation for the Khmer language to thrive in the international digital landscape. With a reliable and consistent data representation, the path is

cleared for the development and advancement of more sophisticated language technologies. Reliable machine translation, accurate text-to-speech and speech-to-text systems, and nuanced natural language processing models all depend on clean, predictable, and canonical text data. By embracing and implementing this standard, the global technology community ensures that the rich literary and cultural heritage of the Khmer script is not only preserved but can become a vibrant and integral part of the future digital world.

## Works cited

1. Khmer script - Wikipedia, accessed on September 1, 2025, [https://en.wikipedia.org/wiki/Khmer\\_script](https://en.wikipedia.org/wiki/Khmer_script)
2. The Origin of the Graph **in the Thai Script**, accessed on September 1, 2025, <http://www.sealang.net/sala/archives/pdf8/ferlus1997origin.pdf>
3. Khmer Script Resources - W3C, accessed on September 1, 2025, <https://www.w3.org/TR/khmr-lreq/>
4. Khmer orthography notes - r12a.io, accessed on September 1, 2025, <https://r12a.github.io/scripts/khmr/km.html>
5. Khmer Script | PDF | Languages Of Asia | Latin Alphabet - Scribd, accessed on September 1, 2025, <https://www.scribd.com/document/160381935/Khmer-Script>
6. Proposal for Khmer Script Root Zone Label Generation Rules (LGR) | icann, accessed on September 1, 2025, <https://www.icann.org/en/system/files/files/proposal-khmer-lgr-15aug16-en.pdf>
7. Khmer (Unicode block) - Wikipedia, accessed on September 1, 2025, [https://en.wikipedia.org/wiki/Khmer\\_\(Unicode\\_block\)](https://en.wikipedia.org/wiki/Khmer_(Unicode_block))
8. “𑄌” U+1780 Khmer Letter Ka Unicode Character - Compart, accessed on September 1, 2025, <https://www.compart.com/en/unicode/U+1780>
9. khmer encoding structure - Institute of Digital Research & Innovation, accessed on September 1, 2025, <https://www.idri.edu.kh/khmer-encoding-structure/>
10. The orthographical syllable in Khmer and rules for the rendering of register shifters - SIL Language Technology, accessed on September 1, 2025, <https://software.sil.org/downloads/r/mondulkiri/Mondulkiri-5.513-Ortho.pdf>
11. Khmer Encoding Structure - Unicode, accessed on September 1, 2025, [https://www.unicode.org/notes/tn61/utn61-Khmer\\_Encoding\\_Structure\\_V2.pdf](https://www.unicode.org/notes/tn61/utn61-Khmer_Encoding_Structure_V2.pdf)
12. UTN #61: Khmer Encoding Structure - Unicode, accessed on September 1, 2025, <https://www.unicode.org/notes/tn61/>
13. Chapter 16 – Unicode 16.0.0, accessed on September 1, 2025, <https://www.unicode.org/versions/Unicode16.0.0/core-spec/chapter-16/>
14. Chapter 1 – Unicode 16.0.0, accessed on September 1, 2025, <https://www.unicode.org/versions/Unicode16.0.0/core-spec/chapter-1/>
15. Announcing The Unicode® Standard, Version 16.0, accessed on September 1, 2025, <http://blog.unicode.org/2024/09/announcing-unicode-standard-version-160.html>
16. page with code points U+1780 to U+17FF - UTF-8 encoding table and Unicode characters, accessed on September 1, 2025, <https://utf8-chartable.de/unicode-utf8-table.pl?start=6016&number=128&utf8=Ox>

17. Unicode Block “Khmer Symbols” - Compart, accessed on September 1, 2025, <https://www.compart.com/en/unicode/block/U+19E0>
18. Khmer Encoding Structure - Unicode, accessed on September 1, 2025, [https://www.unicode.org/notes/tn61/utn61-Khmer\\_Encoding\\_Structure\\_V1.pdf](https://www.unicode.org/notes/tn61/utn61-Khmer_Encoding_Structure_V1.pdf)
19. Unicode Characters in the Khmer Script - UnicodePlus, accessed on September 1, 2025, <https://unicodeplus.com/script/Khmr>
20. Category:Unicode 1780-17FF Khmer - Wikimedia Commons, accessed on September 1, 2025, [https://commons.wikimedia.org/wiki/Category:Unicode\\_1780-17FF\\_Khmer](https://commons.wikimedia.org/wiki/Category:Unicode_1780-17FF_Khmer)
21. Khmer - The Unicode Standard, Version 16.0, accessed on September 1, 2025, <https://www.unicode.org/charts/PDF/U1780.pdf>
22. KhmerU1780.pdf - Khmer fonts, accessed on September 1, 2025, <https://www.khmerfonts.info/bauhahn/KhmerU1780.pdf>
23. www.unicode.org, accessed on September 1, 2025, <https://www.unicode.org/L2/L2022/22290-khmer-encoding.pdf>
24. UAX #29: Unicode Text Segmentation, accessed on September 1, 2025, <https://unicode.org/reports/tr29/>
25. Khmer ordering rules - Open-Std.org, accessed on September 1, 2025, <https://www.open-std.org/jtc1/sc22/wg20/docs/n1076-Khmer-order11.pdf>
26. Khmer Unicode Sorting gamma, accessed on September 1, 2025, <https://www.khmerfonts.info/bauhahn/KhmerSortingUnicodegamma.pdf>