

Discussion points for UNICODE

UTC#185 Meeting

29 October 2025

UNICODE is used for all 22 scheduled languages of India. Data is useful for Mission Bhashini under which various AI based language technologies are being developed by MeitY, GoI. The UNICODE standard is script based, however teachings in schools and implementations by industry are language based.

Background:

India's linguistic and script diversity is rooted in the concept of the **Akshar**, which represents the fundamental orthographic unit in Indic writing systems. Each Akshar typically corresponds to a Grapheme cluster composed of consonants, vowels, and modifiers and functions as the true linguistic building block of Indian languages.

While the **UNICODE** effectively encodes individual characters for all major Indic scripts, it does not explicitly represent the language specific **Akshar-level structure** or the rules that govern how characters combine to form syllabic units. This **character-based encoding** approach limits the faithful digital representation of Indic languages, resulting in:

- Inconsistent rendering and conjunct formation across platforms,
- Difficulties in text search, indexing, and sorting,
- Ambiguities in speech-text alignment and linguistic modeling, and
- Challenges in interoperability across AI, NLP, and transliteration systems.

Purpose of the document:

The purpose of this document is to propose the integration or formal recognition of “*AKSHAR Coded Character Set and Composition Rules Language: Hindi*” within the UNICODE Standard, either as a supplementary guideline or an annexure to the existing Indic script specifications.

The proposal seeks to:

- Introduce standardized rules for Hindi composition rules that complement UNICODE's existing character encoding.
- Provide language-specific framework , currently for Hindi as Akshar Document

Proposed Solution:

To bridge this gap, AKSHAR Coded Character Set and Composition Rules Language: Hindi is a draft standard by the Bureau of Indian Standards (BIS), supported by MeitY and Department of Science and Technology (DST), to define coded character sets and rules for Indian languages starting with Hindi. It establishes formal grammar for valid syllable formation / grapheme cluster rules, ensuring consistent digital text representation across systems. AKSHAR standardizes the handling of special script elements like Nukta, Virama, ZWJ, and ZWNJ, preventing ambiguity in rendering and processing.

By providing validation tools and test resources, it helps developers, font designers, and technology providers ensure interoperability. AKSHAR forms the foundation for digital use of languages using Indian scripts, enabling applications in search, storage, publishing, NLP, and AI/ML. The draft Hindi document has following components:

1. Charset for Hindi language
2. Composition rules as Augmented Backus Naur Formalism (ABNF) for Hindi
3. Sample test strings

Current status:

The Hindi Akshar Document, developed under the supervision of the Ministry of Electronics and Information Technology (MeitY), Department of Science & Technology (DST). It has been formally submitted to the Bureau of Indian Standards (BIS) for review and standardization and available at

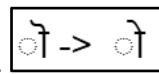
https://www.services.bis.gov.in/tmp/WCLITD23526504_22082025_1.pdf (since archived)

The submission was made under the BIS Technical Committee LITD 20 – Language Technology and Localization, which oversees the development and adoption of standards related to Indian language computing and digital text representation.

Specific Discussion Points

1	A B N F formalism for Devanagari	Existing Unicode Annexure: UAX#29 (Unicode Text Segmentation Report), https://unicode.org/reports/tr29/ Syllable Definition/Grapheme Cluster Rules Useful for Cluster Formation, Rendering, Editing, Cursor / Caret - Movement, Hyphenation / Dropcaps, etc.
---	--	---

2	<p>A B N F formalism for Hindi (Language S p e c i f i c rules)</p>	<p><i>AKSHAR Coded Character Set and Composition Rules Language: Hindi draft</i> (https://www.services.bis.gov.in/tmp/WCLID23526504_22082025_1.pdf)</p> <p>Why Hindi Charset: Devanagari script is shared across languages such as Hindi, Bodo, Dogri, Konkani, Maithili, Marathi, Sanskrit, Santali, Sindhi, etc. Users get confused due to Devanagari Support in Fonts, Keyboards, etc.</p> <div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;"> <p>0912 ओ DEVANAGARI LETTER SHORT O <ul style="list-style-type: none"> • Kashmiri, Bihari languages • also used for transcribing Dravidian short o </p> <p>0913 ओ DEVANAGARI LETTER O</p> </div> <div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;"> <p>094A ॐ DEVANAGARI VOWEL SIGN SHORT O <ul style="list-style-type: none"> • Kashmiri, Bihari languages • also used for transcribing Dravidian short o </p> <p>094B ॐ DEVANAGARI VOWEL SIGN O</p> </div> <p>For example : Users feel that 0912 is artistic rendition of 0913</p> <p>For Hindi Charset, we will include only 0913 and 094B.</p> <p>Or update annotation as below: 0913 DEVANGARI LETTER O: <i>List out languages</i> 094B DEVANGARI VOWEL SIGN O: <i>List out languages</i></p> <p>*Even for single script-language cases, there are archaic characters which are causing confusions</p> <p>Why Hindi Rules:</p> <p>A sequence which is valid in a language, may be invalid in another language. During the process of adding support to more languages, existing implementations may have got broken. Inconsistent implementation across platforms / OS may be due to lack of documentation of language specific behaviour. See reference links below for variations in such rules UNICODE, Font, ICANN, Harfbuzz, W3C, etc. While many of these cater to Devanagari, we have tuned the same for Hindi Character set and Hindi specific behavior including special signs viz. Avagraha and Abbreviation Sign.</p> <p>Examples of language specific rules</p> <div style="display: flex; justify-content: space-around;"> <div style="border: 1px solid black; padding: 10px; text-align: center;"> <p>Examples : Hindi - Santali</p> <p>क् ॐ = क्॑ क्।</p> </div> <div style="border: 1px solid black; padding: 10px; text-align: center;"> <p>Hindi - Kashmiri</p> <p>क् ॐ = क्॑ क्।</p> </div> </div> <p>* Combinations in Red are invalid and Green are valid for respective languages.</p> <p>Ask:</p> <ol style="list-style-type: none"> (1) “AKSHAR Document: Coded Character Set and Composition Rules Language: Hindi , can this document be distributed as UTR/UTN or both? (2) Where in UNICODE can we contribute / enrich such information : UTR / CLDR / Indic mailing list / https://www.unicode.org/ucd/Chapter_12 – Unicode_16.0.0 / UTN (procedure / WG / responsible person) ? (3) UNICODE may guide which WG can we contribute to / who will deal with it (4) Can a focused / dedicated WG be formed so that we can expedite similar exercise for all 22 scheduled languages and encourage rollout of implementations
---	---	---

3	Valid / Invalid Combinations	<p>Visually, the same but encoding different causes problems in Search, TTS output, Sorting.</p> <div style="border: 1px solid black; padding: 10px; margin-bottom: 10px;"> <p>सहयोग = स ह य ंो ग सहयोग = स ह य ंा ंे ग</p> </div> <div style="border: 1px solid black; padding: 10px; margin-bottom: 10px;"> <p>যোগা = য ংো গ ংা যোগা = য ংা ংে গ ংা</p> </div> <p>চাংংদ, কিংতু ব</p> <p>On Notepad in Windows the invalid (2nd string) is rendered with a dotted circle (U+025CC), while in most browsers it joins and both strings look alike.</p> <p>A list for Hindi is submitted for feedback in August/September 2025.</p> <p>https://drive.google.com/file/d/1Sk8mqiUitt0AI2vjK9KERxDjtYPJSTgP/view?usp=drive_link</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="background-color: yellow;">Category</th><th style="background-color: yellow;">For</th><th style="background-color: yellow;">Use</th><th colspan="2" style="background-color: yellow;">Do Not Use</th><th style="background-color: yellow;">Rationale</th></tr> </thead> <tbody> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+0949</td><td>ংং</td><td>U+093E U+0945</td><td>Ambiguity</td></tr> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+0949</td><td>ংংT</td><td>U+0945 U+093E</td><td>Ambiguity</td></tr> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+094A</td><td>ংং</td><td>U+093E U+0946</td><td>Ambiguity</td></tr> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+094A</td><td>ংংT</td><td>U+0946 U+093E</td><td>Ambiguity</td></tr> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+094B</td><td>ংং</td><td>U+093E U+0947</td><td>Ambiguity</td></tr> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+094B</td><td>ংংT</td><td>U+0947 U+093E</td><td>Ambiguity</td></tr> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+094C</td><td>ংং</td><td>U+093E U+0948</td><td>Ambiguity</td></tr> <tr> <td>Indic_Vowel_Sign</td><td>ং</td><td>U+094C</td><td>ংংT</td><td>U+0948 U+093E</td><td>Ambiguity</td></tr> </tbody> </table> <p>Preferred Combinations out of Valid List: Specific cases have also been submitted for consideration in “DoNotEmit” list</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="background-color: yellow;">Category</th><th style="background-color: yellow;">For</th><th style="background-color: yellow;">Use</th><th colspan="2" style="background-color: yellow;">Do Not Use</th><th style="background-color: yellow;">Rationale</th></tr> </thead> <tbody> <tr> <td>Devanagari_Eyelash_Ra</td><td>ং</td><td>U+0931 U+094D</td><td>=</td><td>U+0930 U+094D U+200D</td><td>Compatibility</td></tr> </tbody> </table> <p>Reference Link: http://www.unicode.org/reports/tr39 https://unicode.org/reports/tr29/Combination https://www.unicode.org/Public/17.0.0/ucd/DoNotEmit.txt</p> <p>Ask:</p> <ol style="list-style-type: none"> (1) The submitted list to consider (Status of the list) under DoNotEmit.txt or confusables.txt (Script specific). Any further clarification required to complete the current exercise? (2) Can DoNotEmit be language-wise, so that language specific list can also be proposed <p>E.g.  [DEVANAGARI VOWEL SIGN SHORT O of Kashmiri if typed by user may be normalised to Hindi DEVANAGARI VOWEL SIGN O] These strings can be added in ‘do-not-emit’ for language Hindi.</p> <ol style="list-style-type: none"> (3) How are suggestions verified ? 	Category	For	Use	Do Not Use		Rationale	Indic_Vowel_Sign	ং	U+0949	ংং	U+093E U+0945	Ambiguity	Indic_Vowel_Sign	ং	U+0949	ংংT	U+0945 U+093E	Ambiguity	Indic_Vowel_Sign	ং	U+094A	ংং	U+093E U+0946	Ambiguity	Indic_Vowel_Sign	ং	U+094A	ংংT	U+0946 U+093E	Ambiguity	Indic_Vowel_Sign	ং	U+094B	ংং	U+093E U+0947	Ambiguity	Indic_Vowel_Sign	ং	U+094B	ংংT	U+0947 U+093E	Ambiguity	Indic_Vowel_Sign	ং	U+094C	ংং	U+093E U+0948	Ambiguity	Indic_Vowel_Sign	ং	U+094C	ংংT	U+0948 U+093E	Ambiguity	Category	For	Use	Do Not Use		Rationale	Devanagari_Eyelash_Ra	ং	U+0931 U+094D	=	U+0930 U+094D U+200D	Compatibility
Category	For	Use	Do Not Use		Rationale																																																															
Indic_Vowel_Sign	ং	U+0949	ংং	U+093E U+0945	Ambiguity																																																															
Indic_Vowel_Sign	ং	U+0949	ংংT	U+0945 U+093E	Ambiguity																																																															
Indic_Vowel_Sign	ং	U+094A	ংং	U+093E U+0946	Ambiguity																																																															
Indic_Vowel_Sign	ং	U+094A	ংংT	U+0946 U+093E	Ambiguity																																																															
Indic_Vowel_Sign	ং	U+094B	ংং	U+093E U+0947	Ambiguity																																																															
Indic_Vowel_Sign	ং	U+094B	ংংT	U+0947 U+093E	Ambiguity																																																															
Indic_Vowel_Sign	ং	U+094C	ংং	U+093E U+0948	Ambiguity																																																															
Indic_Vowel_Sign	ং	U+094C	ংংT	U+0948 U+093E	Ambiguity																																																															
Category	For	Use	Do Not Use		Rationale																																																															
Devanagari_Eyelash_Ra	ং	U+0931 U+094D	=	U+0930 U+094D U+200D	Compatibility																																																															

4	CLDR	<p>The Hindi Character set from the Akshar document may be used in CLDR.</p> <p>For language specific Keyboard Layout:</p> <p>https://www.unicode.org/reports/tr35/tr35-keyboards.html https://www.unicode.org/resources/keyboards.html https://cldr.unicode.org/index/keyboard-workgroup</p> <p>Emoji: https://st.unicode.org/cldr-apps/v#/hi/People/</p> <p>Language wise Charset (Currently some characters in link below do not belong to Hindi): https://www.unicode.org/cldr/cldr-aux/charts/28/collation/hi.html</p> <p>Ask:</p> <p>(1) How can we contribute to Keyboard WG? Or attend Keyboard WG meetings?</p> <p>(2) Is the repository of keyboard Public? If not, how can we access and give our feedback?</p> <p>(3) We have an understanding that, there will be rules for accepting keyboards, such as they have to output normalized text and avoid DoNotEmit sequences, otherwise they would not be allowed. If DoNotEmit is script specific, how we can avoid language specific “DoNotEmit” sequences?</p>
8	Inviting voluntary Participation of Unicode in C-DAC Working Groups	<p>We invite experts for voluntary participation in the C-DAC working groups, so that exercise may be taken up for 22 scheduled languages.</p> <p>WG 1: Script Grammar and Script-Language gaps : Language experts + Academia + Implementing agencies</p> <p>WG 2: Hardware & Software (Keyboard, Fonts, Rendering Engine, Applications, Security) : Implementing agencies</p> <p>WG 3: Training & Capacity Building : Awareness and capacity building</p>
9	References	<p>Central Hindi Directorate (CHD) : https://www.chdpublication.education.gov.in/ebook/pdf/devanagarilipiandhindivartanikamankikaran.pdf</p> <p>Other links with grapheme / cluster formation rules</p> <p>BIS - ISCII (IS:13194) Indian Inscript Word syntax : rule 8.1</p> <div data-bbox="466 1311 1019 1558" style="border: 1px solid black; padding: 5px;"> <p>An Indian script word contains one or more syllables, the syntax for which is Backus-Naur Formalism (BNF).</p> <pre> Word ::= {Syllable} [Cons-Syllable] Syllable ::= Cons-Vowel-Syllable Vowel-Syllable Vowel-Syllable ::= V [D] Cons-Vowel-Syllable ::= [Cons-Syllable] Full-Cons [M] [D] Cons-Syllables ::= [Pure-Cons] [Pure-Cons] Pure-Cons Pure-Cons ::= Full-Cons H Full-Cons ::= C [N] </pre> </div> <p>UNICODE Grapheme https://unicode.org/reports/tr29/ (A grapheme is a user-perceived unit of text) and https://www.unicode.org/L2/L2021/21112-deva-cluster-valid.pdf</p>

Expected Outcomes

- **Standardized / consistent Representation:** Language specific Akshar-level encoding rules leading to uniform text behavior across systems.
- **Improved Search and Indexing:** Search engines and NLP pipelines can match, sort, and retrieve words based on Akshar units.
- **Enhanced Linguistic Processing:** Better segmentation and recognition in OCR, ASR, and TTS pipelines for AI engines.
- **Cross-Script Interoperability:** Easier transliteration, translation, and rendering across Indic languages.
- **Dataset Standardization:** Support for consistent corpus creation and linguistic annotation for NLP and AI.

Acknowledgements:

The **Akshar Document** has been developed through extensive consultation and collaborative efforts involving a diverse group of stakeholders including Language Experts, Language Bodies, CDAC, BIS Technical Committee, Academia, Cultural organizations, industry and technology experts.

Way forward:

Submission of the knowledgebase (*AKSHAR Document: Coded Character Set and Composition Rules Language: Hindi*) to the UNICODE Consortium for review and consideration. The intent is to facilitate the reference / integration within the relevant UNICODE Technical Annex (UAX) or associated documentation pertaining to Indic scripts.