

Appendix 3 : Creating the Perso-Arabic tagset

As described in section 2.2.9.2, a tagset using Perso-Arabic characters has been created. The Perso-Arabic tags are given alongside the Roman tags in the descriptions of the tagset in chapter three, and elsewhere the Roman tags alone are used. In this Appendix I outline the principles used to create the Perso-Arabic tags.

Rather than create tags from scratch, I devised a system whereby each character in the Roman tags would map to a given Perso-Arabic character. Creating these mappings as laid out below was by no means unproblematic. To simply map each Roman character to its nearest phonetic equivalent in the Perso-Arabic alphabet would not be useful. Without reference to English grammatical terminology – or the Latin and Greek roots on which much of it is based – the significance of characters mapped in this way would not be clear. Therefore it was necessary to choose characters that “mean something” in the context of Urdu grammatical terminology.

There were two additional difficulties. Firstly, many Urdu grammatical terms begin with the words *harf*, “letter/particle” or *ism*, “name/noun” so that tag initials had to be taken from words modifying these terms. Secondly the grammatical terms I needed and the grammatical terms I was able to find using Haq (2001) did not always coincide. Because of these flaws, this Perso-Arabic version must be considered provisional, pending the approval of some future user whose native language is, if not Urdu, then at least a language written in the Perso-Arabic alphabet.

As a result of this difficulty in choosing appropriate letters, there is a many-to-many mapping between Roman and Perso-Arabic characters, as the table below outlines. Where phonetic / orthographic justification is cited, it means that both the Roman and the Perso-Arabic characters represent some element of an Urdu sound or

letter associated with the decomposable element in question.

Despite the fact that the Perso-Arabic script is cursive, the tags contain only the independent forms of the letters. This is used to parallel the use of capital letters in the Roman tagset. It will also help the reader to distinguish tags visually from the words of the text. In Unicode files, it is to be specified that every character of a Perso-Arabic tag except the last will be followed by the zero-width non-joiner character (U+200C), to ensure that the characters are not joined together when displayed or printed.

Roman	Perso-Arabic	Justification
A (honorific)	ا	Phonetic/orthographic justification
A (in AL)	ا	Phonetic/orthographic justification
A (in AU)	ا	Phonetic/orthographic justification
A (in FA)	ا	I was unable to find an Urdu word for “acronym” and have thus used the transliteration of Roman A.
B (in IB)	ب	Second letter of <i>āgē</i> , “before”
B (in FB)	م	First letter of <i>muxaffaf</i> , “abbreviation”
B (in XB)	ب	Phonetic/orthographic justification
C (cāhiē)	چ	Phonetic/orthographic justification
C(onjunction)	ت	First letter of second word of <i>harfe tardīd</i> , “(disjunctive) conjunction”
C(litic)	چ	First letter of <i>chōTā</i> , “small, short”

C (reciprocal)	پ	Second letter of <i>āpas</i> , the reciprocal pronoun.
D(eterminer)	ت	First letter of second word of <i>harfe ta'rīf</i> , “definite article” ¹
D(egree)	ت	First letter of second word of <i>isme tafzīl</i> , “comparative/superlative adjective”
E	ے	Phonetic/orthographic justification
F(eminine)	ع	First letter of <i>'aurat</i> , “woman”
F(oreign)	د	First letter of <i>dūsrā</i> , “other, alternate”
F(raction)	ب	First letter of <i>batē</i> , “divided”. (The actual word for “fraction” could not be used since its initial letter is already in use in this context.)
F (reflexive)	ج	Second letter of <i>rāja</i> ’, “returning”, part of the Urdu term for “reflexive”
G (for gā)	گی	Phonetic/orthographic justification
G (for gunā)	گی	Phonetic/orthographic justification
G (possessive)	ق	First letter of <i>qābiz</i> , “possessor”
H (for hōnā)	ہ	Phonetic/orthographic justification
H (in XH)	ہ	Phonetic/orthographic justification
I (postposition)	ج	First letter of second word of <i>harfe jār</i> , “preposition”
I (imperative)	ر	Third letter of <i>amro nahī</i> , “imperative mood”

¹ This is not a precise equivalent, but I could not find a better Urdu grammatical term.

J (relative)	ج	Phonetic/orthographic justification
J (adjective)	ص	First letter of second word of <i>isme sifat</i> , “adjective”
K (interrogative)	ک	Phonetic/orthographic justification
L (in AL)	ل	Phonetic/orthographic justification
L (lexical)	ظ	Last letter of <i>lafz</i> , “word”
M (first person)	م	Phonetic/orthographic justification
M(arked)	ا	The letter in which all marked adjectives and many marked nouns end, when in their citation form (it is the masculine nominative inflection)
M(asc)	م	First letter of <i>muzakkar</i> , “male”
M(odal)	ط	First letter of <i>tarīqah</i> , “mode” (probably not in the linguistic sense, but I was unable to find a better word)
N(indefinite)	غ	First letter of <i>ḡairma’ayyan</i> , meaning “indefinite”
N(infinitive)	ن	Phonetic/orthographic justification
N(egative)	ن	Phonetic/orthographic justification (based on the Urdu word for “not”, which is <i>nah</i>).
N(ominative)	خ	First letter of <i>xātī</i> (? ²), “nominative”

² Due to unclear printing in Haq (2001), I was unable to make out the word for “nominative” exactly, but I am reasonably certain that *x* is the first letter.

N(oun)	س	Second letter of <i>ism</i> , “noun”
N(umber)	ع	First letter of <i>'adad</i> , “numeral”
O(blique)	ص	Second letter of <i>nasb</i> , “oblique/accusative case”
O (in FO)	ر	First letter of <i>riyāzī</i> , “mathematics”
O (in OO)		Phonetic/orthographic justification
P(ast)	ض	Third letter of <i>māzī</i> , “past”
P(redicate)	خ	First letter of second word of <i>sifat xabarī</i> , “predicate adjective”
P(ronoun)	ض	First letter of second word of <i>isme zamīr</i> , “pronoun”
P(roper)	ن	First letter of <i>nām</i> , “name”
Q(uestion)	کی	Phonetic/orthographic justification (this is the first letter of all the interrogative words, as well as the question marker)
R (for rahā)	ر	Phonetic/orthographic justification
R (adverb)	ل	Letter occurring twice in <i>muta'aliq fa'l</i> , “adverb” (parallel to R in English, the second consonant in the lexeme “verb”)
R(eflexive or reciprocal)	ر	First letter of <i>rāja</i> ’, “returning”, part of the Urdu term for reflexive
S (for sā)	س	Phonetic/orthographic justification

S(subjunctive)	ش	Third letter of word meaning “subjunctive” ³
S(ubordinating)	ش	First letter of second word of <i>harfe śarT</i> , “conditional conjunction”
S (in FS)	ح	First letter of <i>harf</i> , meaning “letter”
T (2 nd person)	ت	Phonetic/orthographic justification
T (imperfective)	ت	Phonetic/orthographic justification
T (in TT)	ٹ	Second letter of <i>ciT</i> , meaning “tag, label” (from English “chit”)
T (in XT)	ت	Phonetic/orthographic justification
U(nmarked)	ن	Second letter of <i>nah</i> , meaning “not” as in “not marked”
U (in AU)	و	Phonetic/orthographic justification
U (in FU)	ن	First letter of word <i>nāma ’lūm</i> , meaning “unknown”
V (for vālā)	و	Phonetic/orthographic justification
V (3 rd person)	و	Phonetic/orthographic justification
V (far-dem.)	و	Phonetic/orthographic justification
V(erb)	ف	First letter of <i>fa ’l</i> , meaning “verb”
V(ocative)	د	Second letter of second word of <i>harfe nidā</i> , “vocative”

³ Once again, due to unclear printing in Haq (2001), I could not make out the whole word.

X (non-lexical)	غ	Second letter of second word of <i>ke baGair</i> , meaning “without”, as in “without lexical content”. The Urdu grammatical term which is closest to “non-lexical” is probably <i>harf</i> but this term is too general.
X (in FX)	ج	First letter of <i>lafz</i> , “word” (the closest thing I could find to “string” in the computational sense)
Y (near-dem.)	س	Phonetic/orthographic justification
Y (perfective)	س	Phonetic/orthographic justification
Z (in ZZ)	ز	Phonetic/orthographic justification (fullest way to write the izāfat clitic)
Z (in FZ)	ت	First letter of second word of <i>hurūfe tahajjī</i> , “letters of the alphabet”
0 (root form)	۰	The corresponding Urdu numeral is used.
1 (singular)	۱	The corresponding Urdu numeral is used.
2 (plural)	۲	The corresponding Urdu numeral is used.

Punctuation tags are discussed separately in section 3.11.