# CONTEXTUAL SHAPE ANALYSIS OF NASTALIQ

*Aamir Wali, Atif Gulzar, Ayesha Zia, Muhammad Ahmad Ghazali,*

*Muhammad Irfan Rafiq, Muhammad Saqib Niaz, Sara Hussain, and Sheraz*

*Bashir*

## ABSTRACT

Nastaliq calligraphic style is one of the most complex and widely used styles of Urdu script. It employs different shapes and sizes for the same letter in a bewildering variety of contexts. It is observed that these different shapes vary with neighboring letters as well as position of that letter in a ligature. This paper explores this variety in shapes that Nastalique offers, thus providing a foundation to eventually model this inherent variety and recreate this script font as written by a calligrapher.

## 1. INTRODUCTION

Nastaliq is one of the most widely used fonts of Urdu script in Pakistan and is regarded as one of the most complex fonts in the literature of electronic computing. Early efforts to make this font available electronically used the approach of storing all possible ligatures (sequence of Urdu characters occurring without space) of Urdu, which was not supported by the standard Microsoft Windows font technology of TTF (True Type Font) available at that time. Consequently this font was limited to some specialized Urdu Word Processors only. With the emergence of new font technology OTF (Open Type Font) from Microsoft, it is now possible to make Nastaliq font, which would be portable to any standard word processor like Microsoft Word, Power Point and Excel.

A natural approach for this, is to make a character-based Nastaliq font i.e. in spite of storing all possible ligatures, individual characters are stored. Character-based Nastaliq font needs to specify the rules, which govern the change of shape of characters while moving from one segment of ligature to another. The first step of such a complex writing system is to identify the all-possible shapes of characters in ligatures and then study the rules that change the shape in particular context.

The aim of this paper is to identify all possible shapes of characters based upon the context in which they occur in ligatures, thus providing a foundation to eventually model it.

## 2. LITERATURE REVIEW AND PROBLEM STATEMENT

### 2.1 Urdu Writing System

Urdu is a complete language with its own script, which is a mixture of Arabic and Persian script. Urdu script has total 38 alphabets excluding Aerabs (Vowel Marks). Character set and Aerabs of Urdu are shown below in Fig 1.



**Figure 1(a)**        Character Set of Urdu

**Figure 1(b)**      Aerabs of Urdu

From the above figure it can be observed that several of the basic alphabets of Urdu share same shape, these are differentiated only by the placement of dots or diacritic Tuay on the basic shape. This property makes Urdu script fast and easy to write.

Urdu is written in the opposite direction to English i.e. from right to left. An interesting concept about Urdu is that its number system is written form left to right. So Urdu writing system has both the properties of left to right and right to left writing systems as shown in figure 2.
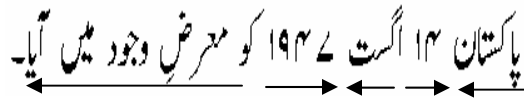
**Figure 2**   Urdu Number System

Urdu writing system is cursive. More than one character joins together to form a ligature. Important thing to be observed is that the characters change their shape depending upon their position in the ligature. Each letter is written in a slightly different form depending on whether it comes in the beginning, middle or end of a word or whether it occurs in isolated form. This is shown in the figure 3.

**Figure 3**     Different shapes according to the position in ligature

In the above figure the character has formed four shapes according to its position.

Another interesting property of Urdu writing system is that characters change their shapes depending upon the characters following and preceding it as in figure 4. This change follows some rules like, shape of the next letter to join with and the shape of the character, which is joining.

**Figure 4**    Different Shapes of the Letter Seen ( س ).

In the above figure it can be clearly seen that the second character ( س ) changes its shape in accordance with following and preceding characters. Which symbolizes that Urdu writing system is context sensitive.

## 2.2 Different Fonts of Urdu

Urdu fonts started growing after the Khat-e-Kofi, an Arabic font that is now being used by Urdu. It was invented in the year 238 AH. But now Kofi has lost its fame because of its complex writing system as shown in figure 5.

**Figure 5**   Different fonts of Urdu

In the year 310 AH a famous Muslim scholar Ibne-Muqalla invented a new and easy font called Naskh. All the previously invented fonts lost their prominence due to Naskh's easy and fast writing system.

Although a considerable number of studies on Urdu Naskh script have been conducted during the last few years and consequently led to the development of software related to Urdu Naskh, studies on Urdu Nastaliq are

still in their premature stages. This is mainly due to extraordinary features of this unique and dynamic form of writing. There are many other fonts of Urdu, some of which are shown in the figure 5.

### 2.3 Urdu Nastaliq Script

Urdu Nastaliq script is a collection of two other scripts Naskh and Talique. These two scripts were combined to form another script formally called Naskh-Talique, which was shortened to Nastaliq. Hence, this 38 alphabets script has properties of both Naskh and Talique.

Two most common feature of Nastaliq found in Naskh or for that matter in any Persian or Arabic script is that it is cursive. That is, the tip of the pen is not raised until a ligature is complete. Another characteristic is that Nastaliq is written from right to left unlike English which is form left to right. In addition to these, there are other characteristics of Nastaliq that have made its automation difficult.

Nastaliq is actually written from top right to bottom left. Each ligature is tilted at approx. 45 degree (See Fig. 6). This is of particular significance as there is no fixed level or height for any character. Positions of characters change at a continuous scale. In Naskh each character has four shapes depending on whether the character is isolated, at start, in middle or coming at end. In Nastaliq, however characters may have significant amount of variation at one position alone. This is because of the influence of other characters in the same ligature.
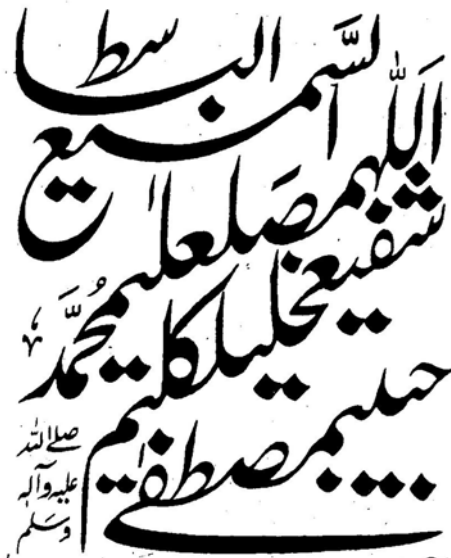


**Figure 6** Nastaliq diagonality

The aim of this paper is to identify all possible shapes of characters of Nastaliq style based upon the context in which they occurs in ligature, thus providing a foundation to eventually model it.

Noori Nastaliq (a type of Nastaliq font, see fig. 5) provides an ideal platform for this kind of analysis. Moving from right to left this style of writing uses simple spacing rules (where size of a shape does not depend on how much space is available for the letter). Thus Noori Nastaliq style of writing eliminates other complexities of Nastaliq font while providing different shapes and their context.

## 3. METHODOLOGY

We need to study all possible ligatures of Noori Nastaliq, to identify different shapes from them. Practically speaking, it is not possible to study all of the ligatures. So we limit our study up to 4 character ligatures (which are about 600,000 ligatures).

We observed that some of them behave similarly under all circumstance, while writing. Which can be observed from the following example in figure 7. In this figure the middle characters have common shapes with the difference of dots.
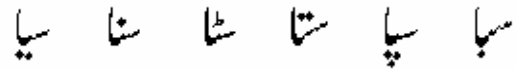


**Figure 7**    Examples of characters behaving similarly in all the contexts

So these 38 characters have been categorized into 21 classes so that any member of same class will have same

shape in ligature provided the same context This minimal subset is shown in Table 1.

**Table 1:** Classification of Urdu characters

ا ـا ـ١

ب پ ت ٹ ث ـ٢

ج چ ح خ ـ٣

د ڈ ذ ـ٤

ر ڑ ز ژ ـ٥

س ش ـ٦

ص ض ـ٧

ط ظ ـ٨

ع غ ـ٩

ف ـ١٠

ق ـ١١

ک گ ـ١٢    ہ ـ١٧

ل ـ١٣    ھ ـ١٨

م ـ١٤    ء ـ١٩

ن ـ١٥    ی ـ٢٠

و ـ١٦    ے ـ٢١

The next step was to identify all possible shapes of all categories. For this purpose a well-known Urdu word processor Inpage version 1.1.34 was used. Typing of large number of combination was a cumbersome and an impossible job. So, a simple program was written that could generate all possible combinations of 2 characters, 3 characters and 4 characters ligature and store them in a text file.

The combinations of 2 to 4 characters were chosen because most of the character in a ligature changes it shape depending upon the shape of the subsequent two to three characters. This can be better understood from the figure 7.



**Figure 7** Change in the shape of a character depending upon the context

In the above figure the character Geem (ج) is in its isolated form. When some other character follows it, it changes its shape, which proves the change in shape according to the next character in 2 characters ligature. Similarly this property can be observed for the 3 and 4 character ligatures.

The text files generated by the program were then imported into the Urdu software so that they can be treated as Urdu characters. We carefully observed all ligatures containing a particular character from the minimal subset that we previously determined (Table 1). Any variation in shape was noted and recorded.

## 4. RESULTS

The following table shows the list of Urdu letters along with the number of shapes that were found for that letter. Appendix A contains the details of the shapes that were found.

The collected data has enabled us to carry another study to investigate the rules that might have caused the change in shape of these characters. This will help us to better understand the nature of Nastaliq.

## 7. REFERENCES

Ethnologue: Languages of the World, 14th Edition, 2001 © Summer Institute of Linguistics, URL: I http://www.sil.org

| No. of Shapes | Letter | No. of Shapes | Letter |
|---|---|---|---|
| 64 | ق | 2 | ا |
| 44 | ک | 52 | ب |
| 29 | ل | 37 | ج |
| 31 | م | 3 | د |
| 52 | ن | 3 | ر |
| 1 | و | 33 | س |
| 34 | ہ | 37 | ص |
| 28 | ﷲ | 22 | ط |
| 52 | ی | 50 | ع |
| 52 | ے | 64 | ف |

## 5. DISCUSSION

One major factor that affects the accuracy of above results is the technique of 'Perceptual examination' i.e. the decision that two shapes in different ligatures are same or not, are identified by human perception.

The other source of inaccuracy could be the fact that Inpage has stored all the ligatures of Noori Nastaliq written by hand. So we might have identified some additional shapes caused by the variation of hand written.

## Appendix A

اردو نستعلیق کے حروفِ تہجی کی تمام ممکن اشکال

| | | |
|---|---|---|
| ا | : | 2 |
| ب | : | 52 |
| ج | : | 37 |
| و | : | 3 |
| ر | : | 3 |
| س | : | 33 |
| ص | : | 37 |
| ط | : | 22 |
| ع | : | 50 |
| ف | : | 64 |
| ق | : | *1 |
| ک | : | 44 |
| ل | : | 29 |
| م | : | 31ا |
| ن | : | *1 |
| و | : | 1 |
| ہ | : | 34 |
| ھ | : | 28 |
| ی | : | *1 |
| ے | : | *1 |

*ب، ن، ے، ی کی اور میان اور شروع میں ایک ہی شکل ہوتی ہے۔ اسی طرح ف، ق کی بھی ایک ہی شکل ہوتی ہے۔

**Total : 474**

# ا

| | | ۱ـ ا | ۲ـ با |

# ب

| ۱ـ با | ۲ـ بب | ۳ـ بج | ۴ـ بر | ۵ـ بس |
| ۶ـ بص | ۷ـ بط | ۸ـ بع | ۹ـ بف | ۱۰ـ بق |
| ۱۱ـ بم | ۱۲ـ بن | ۱۳ـ بھ | ۱۴ـ بہ | ۱۵ـ بی |
| ۱۶ـ بے | ۱۷ـ ببا | ۱۸ـ ببر | ۱۹ـ ببن | ۲۰ـ ببے |
| ۲۱ـ ببع | ۲۲ـ ببق | ۲۳ـ ببج | ۲۴ـ ببکب | ۲۵ـ ببہا |
| ۲۶ـ ببب | ۲۷ـ ببک | ۲۸ـ ببیب | ۲۹ـ ببج | ۳۰ـ ببر |
| ۳۱ـ ببہ | ۳۲ـ ببی | ۳۳ـ ببے | ۳۴ـ ببس | ۳۵ـ ببیق |
| ۳۶ـ ببکا | ۳۷ـ ببکب | ۳۸ـ ببند | ۳۹ـ ببنبل | ۴۰ـ ببیب |
| ۴۱ـ ببیبہ | ۴۲ـ ببقا | ۴۳ـ ببقہ | ۴۴ـ ببیہ | ۴۵ـ ببنک |
| ۴۶ـ ببنا | ۴۷ـ ببنی | ۴۸ـ ببنے | ۴۹ـ ببنج | ۵۰ـ ببسبس |
| ۵۱ـ ببلبد | ۵۲ـ ببطبم | | | |

# ج

| | | | | |
|---|---|---|---|---|
| ۵۔ جس | ۴۔ جر | ۳۔ جچ | ۲۔ جب | ۱۔ جا |
| ۱۰۔ جن | ۹۔ جم | ۸۔ جف | ۷۔ جج | ۶۔ جس |
| ۱۵۔ جبی | ۱۴۔ جبک | ۱۳۔ جبر | ۱۲۔ جج | ۱۱۔ جو |
| ۲۰۔ جها | ۱۹۔ جکل | ۱۸۔ جکب | ۱۷۔ جفا | ۱۶۔ جسے |
| ۲۵۔ جنجا | ۲۴۔ جیل | ۲۳۔ جیب | ۲۲۔ جها | ۲۱۔ جهب |
| ۳۰۔ بجر | ۲۹۔ بج | ۲۸۔ بجب | ۲۷۔ بجا | ۲۶۔ جچ |
| ۳۵۔ بجا | ۳۴۔ بجہ | ۳۳۔ بجج | ۳۲۔ بجس | ۳۱۔ بجس |
| | | ۳۷۔ بجکب | ۳۶۔ بجبر | |

# د

| | | |
|---|---|---|
| ۳۔ سد | ۲۔ بد | ۱۔ د |

# ر

| | | |
|---|---|---|
| ۳۔ سر | ۲۔ بر | ۱۔ ر |

# س

| ۵ـ سس | ۴ـ سر | ۳ـ سج | ۲ـ سب | ۱ـ سا |
|---|---|---|---|---|
| ۱۰ـ سن | ۹ـ سم | ۸ـ سع | ۷ـ سط | ۶ـ سص |
| ۱۵ـ سبب | ۱۴ـ سے | ۱۳ـ سی | ۱۲ـ سہ | ۱۱ـ سو |
| ۲۰ـ سعف | ۱۹ـ سطی | ۱۸ـ سبد | ۱۷ـ سج | ۱۶ـ سمن |
| ۲۵ـ سہ | ۲۴ـ سنب | ۲۳ـ سبر | ۲۲ـ سلک | ۲۱ـ سکب |
| ۳۰ـ سلسم | ۲۹ـ سجج | ۲۸ـ بسو | ۲۷ـ سہس | ۲۶ـ نسفق |
| | | ۳۳ـ سمنس | ۳۲ـ یمسر | ۳۱ـ سہا |

# ص

| ۵ـ صر | ۴ـ صد | ۳ـ ص | ۲ـ صب | ۱ـ صا |
|---|---|---|---|---|
| ۱۰ـ صبا | ۹ـ ص | ۸ـ صو | ۸ـ صم | ۶ـ صق |
| ۱۵ـ صصل | ۱۴ـ صص | ۱۳ـ صبی | ۱۲ـ صبس | ۱۱ـ صبر |
| ۲۰ـ صکا | ۱۹ـ صعب | ۱۸ـ صعا | ۱۷ـ صطر | ۱۶ـ صصا |
| ۲۵ـ صید | ۲۴ـ صہا | ۲۳ـ صہا | ۲۲ـ صنے | ۲۱ـ صلب |

۲۶۔بصا    ۲۷۔یصب    ۲۸۔بصج    ۲۹۔بصر    ۳۰۔صم

۳۱۔بصو    ۳۲۔بصہ    ۳۳۔بصی    ۳۴۔بصھ    ۳۵۔سو

۳۶۔فصو    ۳۷۔لصّف

# ط

۱۔طا    ۲۔طب    ۳۔طج    ۴۔طر    ۵۔طس

۶۔طص    ۷۔طع    ۸۔طف    ۹۔طق    ۱۰۔طم

۱۱۔طہ    ۱۲۔طی    ۱۳۔طے    ۱۴۔طبا    ۱۵۔طبر

۱۶۔طبق    ۱۷۔طبو    ۱۸۔طبہ    ۱۹۔طبی    ۲۰۔طمب

۲۱۔طہا    ۲۲۔بط

# ع

۱۔عا    ۲۔عب    ۳۔عج    ۴۔عر    ۵۔عس

۶۔عص    ۷۔عع    ۸۔عم    ۹۔عو    ۱۰۔عی

۱۱۔عے    ۱۲۔عبا    ۱۳۔عبج    ۱۴۔عبد    ۱۵۔عبر

۱۶۔عجب    ۱۷۔عصا    ۱۸۔عفا    ۱۹۔بع    ۲۰۔بجا

۲۱ـبُجّ     ۲۲ـبعس     ۲۳ـبعہ     ۲۴ـبعے     ۲۵ـجعد

۲۶ـجعو     ۲۷ـسعو     ۲۸ـصعب     ۲۹ـصعر     ۳۰ـصعو

۳۱ـطعن     ۳۲ـفعا     ۳۳ـمعر     ۳۴ـنعو     ۳۵ـبعبر

۳۶ـبعکب     ۳۷ـسعید     ۳۸ـطعنو     ۳۹ـطعنے     ۴۰ـکعب

۴۱ـکعبے     ۴۲ـلعبا     ۴۳ـلعفو     ۴۴ـلعیب     ۴۵ـمعبو

۴۶ـمعجم     ۴۷ـمعمو     ۴۸ـنعبد     ۴۹ـنعلو     ۵۰ـبعجی

# ف

۱ـفا     ۲ـفب     ۳ـفج     ۴ـفد     ۵ـفر

۶ـفص     ۷ـفس     ۸ـفع     ۹ـفف     ۱۰ـفق

۱۱ـفک     ۱۲ـفل     ۱۳ـفم     ۱۴ـفن     ۱۵ـفو

۱۶ـفی     ۱۷ـفے     ۱۸ـفہ     ۱۹ـفبب     ۲۰ـفبر

۲۱ـفبس     ۲۲ـفجو     ۲۳ـفجب     ۲۴ـفجم     ۲۵ـفجر

۲۶ـفجھ     ۲۷ـفسا     ۲۸ـفصو     ۲۹ـفصے     ۳۰ـفصا

۳۱ـفطر     ۳۲ـفعب     ۳۳ـفکا     ۳۴ـفکل     ۳۵ـفکس

۳۶ـفلف     ۳۷ـفبا     ۳۸ـفم     ۳۹ـفنا     ۴۰ـفنر

۴۵ـ بفا | ۴۴ـ لبف | ۴۳ـ فیف | ۴۲ـ فہا | ۴۱ـ فہص

۵۰ـ بفق | ۴۹ـ بفص | ۴۸ـ بفس | ۴۷ـ بفر | ۴۶ـ بفج

۵۵ـ بفبس | ۵۴ـ بفبر | ۵۳ـ بفبج | ۵۲ـ بفبب | ۵۱ـ بفبا

۶۰ـ بفبج | ۵۹ـ بفبب | ۵۸ـ بفبا | ۵۷ـ بفبے | ۵۶ـ بفبو

| | | ۶۴ـ جفنا | ۶۳ـ بفبھ | ۶۲ـ بفکا | ۶۱ـ بفبھ

# ک

۵ـ کر | ۴ـ کد | ۳ـ کج | ۲ـ کب | ۱ـ کا

۱۰ـ کم | ۹ـ کک | ۸ـ کع | ۷ـ کص | ۶ـ کس

۱۵ـ کے | ۱۴ـ کھ | ۱۳ـ کہ | ۱۲ـ کو | ۱۱ـ کن

۲۰ـ کبے | ۱۹ـ کبی | ۱۸ـ کبج | ۱۷ـ کبس | ۱۶ـ کبر

۲۵ـ کس | ۲۴ـ کج | ۲۳ـ کجک | ۲۲ـ کجب | ۲۱ـ کجد

۳۰ـ کا | ۲۹ـ بفقا | ۲۸ـ کفا | ۲۷ـ کعر | ۲۶ـ کعج

۳۵ـ کنہ | ۳۴ـ کمب | ۳۳ـ ککو | ۳۲ـ کگر | ۳۱ـ ککے

۴۰ـ کعج | ۳۹ـ کھل | ۳۸ـ کھا | ۳۷ـ کہب | ۳۶ـ کہا

| | | ۴۴ـ بنمک | ۴۳ـ بیاک | ۴۲ـ کا | ۴۱ـ کھص

# ل

| ۵ـلد | ۴ـلج | ۳ـلب | ۲ـلا | ۱ـل |
|---|---|---|---|---|
| ۱۰ـلح | ۹ـلط | ۸ـلص | ۷ـلس | ۶ـلر |
| ۱۵ـلم | ۱۴ـلل | ۱۳ـلک | ۱۲ـلق | ۱۱ـلف |
| ۲۰ـلجد | ۱۹ـلی | ۱۸ـله | ۱۷ـلھ | ۱۶ـلن |
| ۲۵ـللہ | ۲۴ـبلب | ۲۳ـللہ | ۲۲ـلم | ۲۱ـلس |
| | ۲۹ـبلھب (ٮ) | ۲۸ـبلھ | ۲۷ـبلب | ۲۶ـبلکل |

(ٮ اس شکل میں لل چھوٹی ہوتی ہے۔)

# م

| ۵ـمس | ۴ـمر | ۳ـمج | ۲ـمب | ۱ـما |
|---|---|---|---|---|
| ۱۰ـمی | ۹ـمہ | ۸ـمن | ۷ـمع | ۶ـمص |
| ۱۵ـمصا | ۱۴ـمسی | ۱۳ـمج | ۱۲ـمبد | ۱۱ـمے |
| ۲۰ـمما | ۱۹ـملب | ۱۸ـمکب | ۱۷ـمطر | ۱۶ـمطا |
| ۲۵ـجمی | ۲۴ـجمو | ۲۳ـجمن | ۲۲ـجمع | ۲۱ـم |

٣٠۔ سمن    ٢٩۔ سمس    ٢٨۔ سما    ٢٧۔ جبر    ٢٦۔ جے

٣١۔ ہمم

و

ا۔ و

ہ

٥۔ ہر    ٤۔ ہج    ٣۔ ہب    ٢۔ ہا    ١۔ ہ

١٠۔ ہن    ٩۔ ہف    ٨۔ ہج    ٧۔ ہص    ٦۔ ہس

١٥۔ ہجد    ١٤۔ ہبد    ١٣۔ ہج    ١٢۔ ہہ    ١١۔ ہھ

٢٠۔ ہہب    ١٩۔ ہہا    ١٨۔ ہہا    ١٧۔ ہ کل    ١٦۔ ہکا

٢٥۔ ہج    ٢٤۔ ہیص    ٢٣۔ ہیس    ٢٢۔ ہہر    ٢١۔ جج

٣٠۔ ہہو    ٢٩۔ ہہن    ٢٨۔ ہہم    ٢٧۔ ہہق    ٢٦۔ ہہف

٣٤۔ اللہ    ٣٣۔ ہہ    ٣٢۔ ہہہ    ٣١۔ ہھھ

ھ

| | | | | |
|---|---|---|---|---|
| ۵_ھس | ۴_ھر | ۳_ھج | ۲_ھب | ا_ھا |
| ۱۰_ھو | ۹_ھن | ۸_ھم | ۷_ھق | ۶_ھص |
| ۱۵_ھبا | ۱۴_ھے | ۱۳_ھی | ۱۲_ھھ | ۱۱_ھہ |
| ۲۰_ھبے | ۱۹_ھبو | ۱۸_ھبر | ۱۷_ھبج | ۱۶_ھبب |
| ۲۵_ھا | ۲۴_ھکل | ۲۳_ھفھ ۲۲_ھجب | | ۲۱_ھجا |
| | | ۲۸_ھپھ ۲۷_ھباھ | | ۲۶_ھہا |