

# Urdu Morphology, Orthography and Lexicon Extraction

اردو

Presented by:

**Muhammad Humayoun**

Department of Mathematics  
University of Savoie  
France  
*mhuma@univ-savoie.fr*

Co-authors:

**Harald Hammarström**  
**Aarne Ranta**

Department of Computer Science  
Chalmers University of Technology  
Sweden  
*{harald2, aarne}@cs.chalmers.se*

LAMA



UNIVERSITE  
CHAMBERY ANNECY DE SAVOIE

CAASL-2, Stanford



# Introduction

- ***Indo-European → Indo-Iranian → Indo-Aryan***
- Written from **right to left** using **Perso-Arabic** Script.
- **Grammar** and **Vocabulary** influenced by Arabic, Persian and the native languages of South Asia
- Widely Spoken in Pakistan, India and Jammu & Kashmir
- Also spoken all over the world due to big south Asian Diaspora
- **Urdu-Hindi**: share grammar, almost all phonology and lot of vocabulary
- Urdu-Hindi together is the **second most widely spoken language**  
(Native + second language)

# Contribution

- **Orthography component:** A Unicode Infrastructure to accommodate Perso-Arabic script of Urdu
- **Morphology component :**
  - A **type system** that covers the language abstraction completely
  - An **inflection engine** that covers word-and-paradigm morphological rules for all word classes
- **Lexicon:** Automatically extracted, 4,816 words generating 137,182 word forms.
- **Grammar component :** A small fragment of syntax

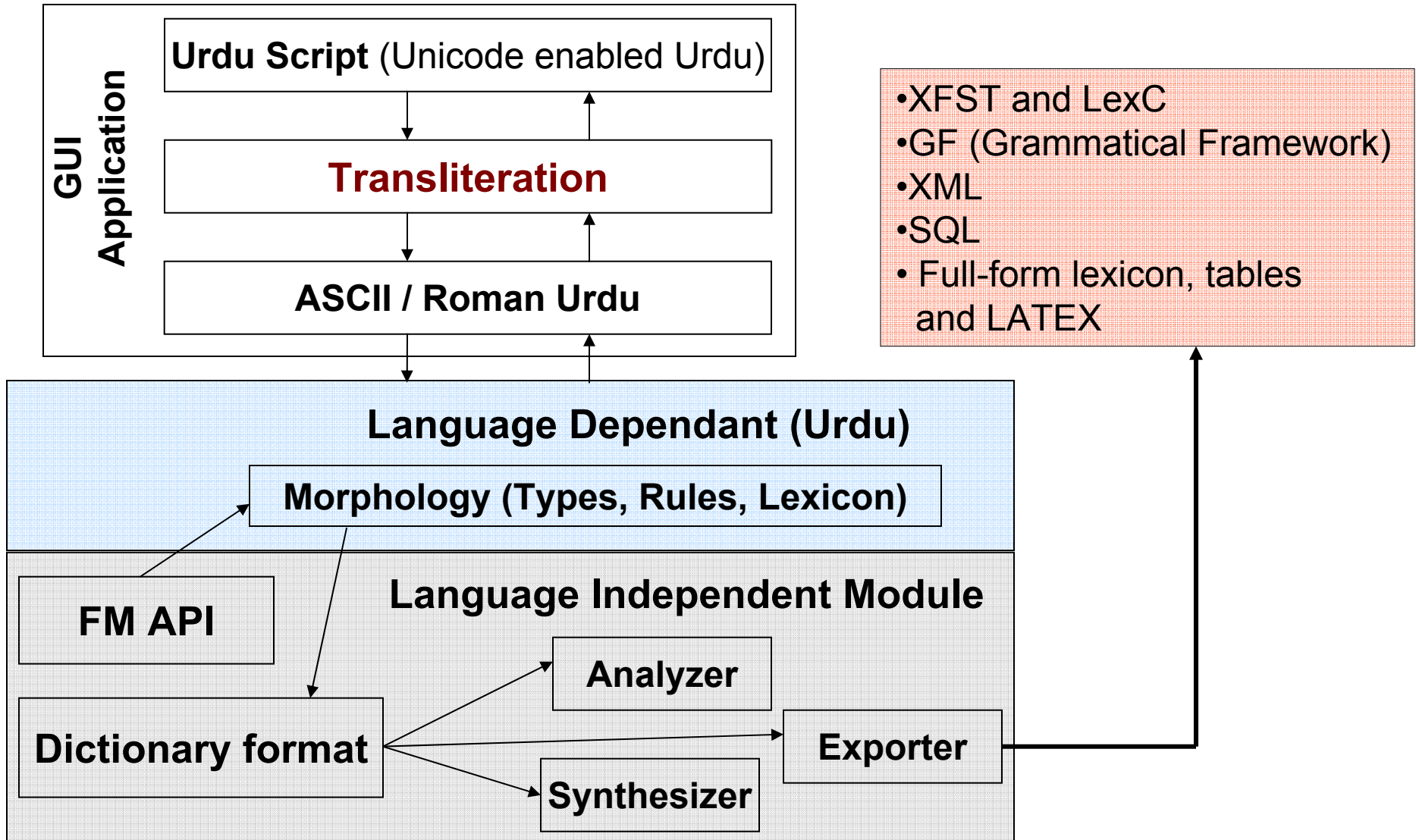
# Urdu Orthography

- An alphabet of 57 letters and 15 diacritic marks
- The use of diacritic marks: optional
- **Morphology** and the **lexicon** saved in ASCII characters
  - **Reusability** for Hindi in future, by adding a lexicon and the transliteration scheme
  - Easy manipulation on different platforms
- **Unicode support** provided by a clear, strict and reversible transliteration scheme (Transliterator)
- A GUI application and useful tools
  - (Keyboard input method, Urdu Extractor)
- Implemented in Java by using ICU4J and Swing packages

# Urdu Morphology

- Morphology is implemented in **Functional Morphology (FM)**
- An open source toolkit or domain embedded language for morphology development in Haskell
  - Functional Programming language, High level of abstraction, Higher-order functions, type classes, polymorphism
  - These features: good for capturing linguistic generalizations
- Idea: Dealing with grammars as reusable software libraries
- Functional Morphology treats
  - The part of speech (word classes) as **data types**
  - Their Inflection as **finite functions**
- Tools (API functions, Analyzer, Synthesizer, Exporter)

# Morphology + Orthography



**Functional Morphology Toolkit**

# Nouns in Urdu: type system

- Urdu Noun Inflects in **Number** (*Singular, Plural*) and **Case** (*Nominative, Oblique, Vocative*)

```
data Number      = Singular | Plural
                  deriving (Show, Eq, Enum, Ord, Bounded)
```

```
data Case        = Nominative | Oblique | Vocative
```

```
....
```

```
data NounForm   = NF Number Case
```

```
....
```

```
type Noun       = NounForm → Str
```

```
....
```

- Inherent parameter: **Gender** (*Masculine, Feminine*)

```
data Gender     = Masculine | Feminine
                  deriving (Show,Eq,Enum,Ord,Bounded)
```

# Nouns in Urdu

- Nouns are divided into **15 groups** based on their inflection
- **A group as running example:**
  - Singular **masculine** nouns ending with ( ل, a) , ( ہ, h) and ( ع, e / ? / ə)
- **Making:**
  - If a word ends with letter ( ل, a) or ( ہ, h) then:
    - **Plural nominative, singular oblique:** last letter is replaced by ( ے, e)
    - **Plural oblique:** the last letter is replaced by ( وں, on)
    - **Plural vocative:** last letter is replaced by ( و, o)
  - If a word ends with ( ع, e / ? / ə): above mentioned letters will be added at the end without replacing any existing letter



# Nouns in Urdu: inflection engine

## Example Noun: (laṛka, لڑکا, boy)

	<i>Nominative</i>	<i>Oblique</i>	<i>Vocative</i>
<i>Singular</i>	laṛka لڑکا	laṛke لڑکے	laṛke لڑکے
<i>Plural</i>	laṛke لڑکے	laṛkon لڑکوں	laṛko لڑکو

noun\_IRka :: DictForm → Noun

noun\_IRka IRka nf = **mkNoun** sg pl sg\_obl pl\_obl sg\_voc pl\_voc nf

where

sg = IRka

pl = IRk ++ "E"

sg\_obl = pl

pl\_obl = IRk ++ "wN"

sg\_voc = pl

pl\_voc = IRk ++ "w"

IRk = if (end == "e") then IRka else (tk 1 IRka)

end = dp 1 IRka

# Nouns in Urdu: inflection engine

- An general function for the Inflection table of nouns

**mkNoun**::String→String→String→String→String→String→Number→Case → String

**mkNoun** sg pl sg\_Obl pl\_Obl sg\_Voc pl\_Voc n c =

case n of

Singular → case c of

Nominative → sg

Oblique → sg\_Obl

Vocative → sg\_Voc

Plural → case c of

Nominative → pl

Oblique → pl\_Obl

Vocative → pl\_Voc

	<i>Nom</i>	<i>Obl</i>	<i>Voc</i>
<i>Sg</i>	laɾka لڑکا	laɾke لڑکے	laɾke لڑکے
<i>Pl</i>	laɾke لڑکے	laɾkon لڑکوں	laɾko لڑکو

# Nouns in Urdu: inflection engine

- An interface function for this group of nouns

n1 :: DictForm → Entry

n1 df = masculine (noun\_IRka df)

*DictForm*: a string type

*masculine*: a function for masculine words

- Defined in the Lexicon:

**n1 l(a)R'ka** (ləɾka, لَرْڪَا)

laɾka  
N  
Masc  
NF Sg Nom: laɾka  
NF Sg Obl: laɾke  
NF Sg Voc: laɾke  
NF Pl Nom: laɾke  
NF Pl Obl: laɾkon  
NF Pl Voc: laɾko

# Nouns in Urdu: inflection engine

- An interface function for this group of nouns

n1 :: DictForm → Entry

n1 df = masculine (noun\_IRka df)

*DictForm*: a string type

*masculine*: a function for masculine words

- Defined in the Lexicon:

**n1 l(a)R'ka** (ləɾka, لَرَّڪَا)

## Some other noun groups in the Lexicon:

**n2 k(o)n'waN** (Well, kuŋwan, گُنْوَان)

**n3 m(a)r'd** (Man, mərd, مَرْد)

**n4 k(o)r'sy** (Chair, kuɾsi, كَرْسِي)

**n5 maN** (Mother, maŋ, مَان)

**n6 g(o)R'ya** (Doll, guɾiya, گُرِّيَا)

**n7 KwX'bw** (Fragrance, xuʃbu, خَوْشْبُو)

**n8 k(i)tab** (book, kiɾab, كِتَاب)

.....

# Urdu Verbs

	Root	Infinitive	Oblique
<b>Intransitive</b> / Transitive / Ditransitive etc	bən بن	<b>bəna</b> بننا <i>to build (by unknown)</i>	bəne بننے
Direct Causative	bəna بنا	<b>bəna</b> بنانا <i>to build (by self)</i>	bəne بنانے
Indirect Causative	bəna بناوا	<b>bəna</b> بنوانا <i>to build (by third person)</i>	bəne بنوانے

We divide verbs in the following categories:

- Basic stem form, direct & indirect causatives exist
- Only Basic stem form exists
- Basic stem form & direct causative form exist
- Basis stem form & indirect causative form exist
- **6 groups** have been implemented for verbs

# Urdu Verbs

- Urdu verb inflects in:
  - Gender, Number
  - Person (First, Second {*casual, familiar, respectful*}, Third {*near, distant*})
  - Tense (Subjunctive, Perfective, Imperfective)

# Urdu Verbs: type system

- Category: *Basic stem form, direct & indirect causatives exist*

```
type Verb = VerbForm → Str
```

```
data VerbForm =
```

```
  VF Tense Person Number Gender |  
  Caus1 Tense Person Number Gender |  
  Caus2 Tense Person Number Gender |  
  Inf          | Caus1_Inf          | Caus2_Inf          |  
  Inf_Fem     | Caus1_Inf_Fem | Caus2_Inf_Fem |  
  Inf_Obl     | Caus1_Inf_Obl | Caus2_Inf_Obl |  
  Root        | Caus1_Root    | Caus2_Root
```

```
data Person = Pers1 |
```

```
  Pers2_Casual | Pers2_Familiar | Pers2_Respect |  
  Pers3_Near | Pers3_Distant
```

```
data Tense = Subj | Perf | Imperf
```

# Urdu Verbs: inflection engine

```
mkVerbCaus12 :: String -> String -> String -> Verb
mkVerbCaus12 vInf caus1_inf caus2_inf =
  mkGenVerb root r1 r2 vInf caus1_inf caus2_inf
  where
    root      = (tk 2 vInf)
    r1        = (tk 2 caus1_inf)
    r2        = (tk 2 caus2_inf)
```

## ■ An general function for the Inflection table

```
mkGenVerb::DictForm → DictForm → DictForm → DictForm → DictForm → DictForm → Verb
mkGenVerb root r1 r2 vf caus1 caus2 (Root1) = root
mkGenVerb root r1 r2 vf caus1 caus2 (VF1 t p n g) = mkVAnalysis root t p n g
```

```
mkVAnalysis :: String → Tense → Person → Number → Gender → String
mkVAnalysis root tense p n g =
  case tense of
    Subjunctive → case p of
      Pers1 → case n of
        Singular → case g of
          Masculine → mkEnding b root "w^N" "wN"
    where
      t = dp 1 root
      b = inStr t ["A","a","w"]
```



# Urdu Verbs: in lexicon

## ■ v4 **bna bnana bnwana**

( bəna بنا , bəna bnana بنانا , bənwana بنوانا )

Root1: **bən** بن

Inf1: **bəna** بننا

Inf\_Obl1: **bəne** بننے

Inf\_Fem1: **bənni** بننی

Caus1\_Root: **bəna** بنا

Caus1\_Inf: **bəna bnana** بنانا

Caus1\_Inf\_Obl: **bəna ne** بنانے

Caus2\_Root: **bənwa** بنوا

Caus2\_Inf: **bənwana** بنوانا

Caus2\_Inf\_Obl: **bənwane** بنوانے

VF Subj Pers1 Sg Masc/ Fem : **bənʊ** بنوں

VF Subj Pers1 Pl Masc/ Fem : **bənən** بنیں

VF Subj Pers2\_Casual Sg Masc/ Fem : **bən** بن

VF Subj Pers2\_Casual Pl Masc/ Fem : **bənʊ** بنو

....

Caus1 Subj Pers1 Sg Masc/Fem: **bənaʔʊ** بناؤں

Caus1 Subj Pers1 Pl Masc/Fem: **bənaʔən** بنائیں

Caus1 Subj Pers2\_Casual Sg Masc/Fem: **bəna** بنا

Caus1 Subj Pers2\_Casual Pl Masc/Fem: **bənaʔo** بناؤ

.....

Caus2 Subj Pers1 Sg Masc/Fem: **bənwaʔʊ** بناؤں

Caus2 Subj Pers1 Pl Masc/Fem: **bənwaʔən** بنوائیں

Caus2 Subj Pers2\_Casual Sg Masc/Fem: **bənwa** بنوا

Caus2 Subj Pers2\_Casual Pl Masc/Fem: **bənwaʔo** بناؤ

.....

# Other word classes

- Adjectives
- Adverbs
- The closed classes
  - Pronouns, PostPositions, Particles, Interjunctions, Conjunctions, Negations, Questions and Numerals

# The Lexicon

- A wide-coverage lexicon is a key part of any morphological implementation
- Aim: to build a lexicon automatically with minimal human efforts
- A tool ***extract*** is used which is provided with the Functional Morphology
- It requires a ***paradigm file*** and a ***corpus***
- To build a ***corpus***:
  - A reasonable amount of Urdu Unicode text was collected from the web (news and literature domain)
  - All the html tags & other non-related information were thrown away by a tool (developed with this work) and save the file as text file
  - Urdu Unicode text is then converted into ***ASCII Urdu*** by transliteration tool
- ***The lexicon***: extracted by applying *paradigms* on *corpus*

# The Lexicon: Problems

- Urdu is commonly written without or with a variant number of diacritic marks
  - A fundamental limitation to get a fully vocalized corpus
- Problem: having **more versions per word** with different diacritics
  - e.g. (كِتَاب, kīṭab) and (کتاب, kṭab) for word (kīṭab, book)
  - **Point:** We should save only one version per word with full diacritics
- Tokens with different diacritics are **not always same words**
  - e.g. (تیر, ṭīr, to swim) and (تیر, ṭīr, arrow)
  - **Point :** We should save all such words with full diacritics

# The Lexicon Extraction - Results

- To assure the correctness:
  - Manually re-checking of the lexicon from word to word
  - Incorrect entries thrown away
- A fundamental limitation
  - The missing diacritics on partly vocalized words are not applied

Corpus	
Size (Words)	1,520,000 (1.5 million)
Words containing Diacritics:	23,696
Unique tokens:	63,700 (4.1%) **

Lexicon	
Extracted lexicon	9,126
Words containing Diacritics:	632
Clean lexicon	4,816 (52.8%)
Words containing Diacritics:	415

\*\*This conforms well to our intuition that **high frequent items** (postpositions, auxiliaries, particles and pronouns), account for most tokens in Urdu text.

# The Lexicon Extraction - Results

- Why so many incorrect entries?
- The strictness of rules in paradigm file: normal
  - Trade-off: quality vs. coverage
- Spelling mistakes:
  - Original Typos
  - Lack of spaces between words
  - Extra spaces inside words
  - **Possible Reason**: The use of Urdu on web is relatively new
- Foreign words:
  - Arabic – The verses of Holy Quran in religious text
  - Persian – Poetry in slightly old literary text
  - Lot of **proper nouns** and **English words** in the news domain

# The Syntax

- Urdu an **SOV** (Subject Object Verb) language
- Relatively **free word order**
- A small fragment of syntax as a separate component on top of morphology by using Grammatical Framework
- Grammatical Framework (GF):
  - A logical Framework
  - Programming language for defining grammars (formal + natural)
  - Grammar = The Abstract syntax and Concrete syntax
- In our Implementation: A sentence:
  - Combination a noun phrase (NP) and a verb phrase (VP)
  - Combination of two sentences by adding a conjunction in between

# The Syntax

## ■ Abstract Syntax:

fun UsePresS: NP → VP → S;

## ■ Concrete Syntax:

UsePresS np vp =

{s = np.s ! Nom ++ vp.s ! Present ! np.p ! np.n ! np.g}

is ko kṛāḅeṇ leni heṇ, اس کو کتابیں لینی ہیں, He/she is suppose to take the books

**DemPron → Num → CN → NP**

ye do kṛāḅeṇ, یہ دو کتابیں, these two books

**DemPron → PN → NP**

wo Ali, وہ علی, that Ali

**NP → PostP → CN → NP**

is ko kṛāḅeṇ, اس کو کتابیں, to him the books

**Verb\_Aux → VP**

heṇ, ہیں, are

**Verb → Verb\_Aux → VP**

leni thiṇ, لینی تھیں, was suppose to take



# A Complete Example

is ko kṭābenḡ leni hen) اِس کو کتابیں لینی ہیں

**Transliteration:** a(i)s kw ktabyN lyny hyN

< اِس, **a(i)s** >

yih\_6 یہ +DemPron - Sg Obl - Pers3\_Near

mayN\_8 میں +PersPron - Sg Pers3\_Near Obl

< کو, **kw** >

kw\_18 کو +PostP

< کتابیں, **ktabyN** >

ktab\_824 کتاب +N - Pl Nom - Fem

< لینی, **lyny** >

lyna\_2 لینا +Verb - Inf\_Fem

< ہیں, **hyN** >

hwana\_0 ہونا +Verb\_Aux - Present Pers1 Pl Masc

hwana\_0 ہونا +Verb\_Aux - Present Pers1 Pl Fem

...

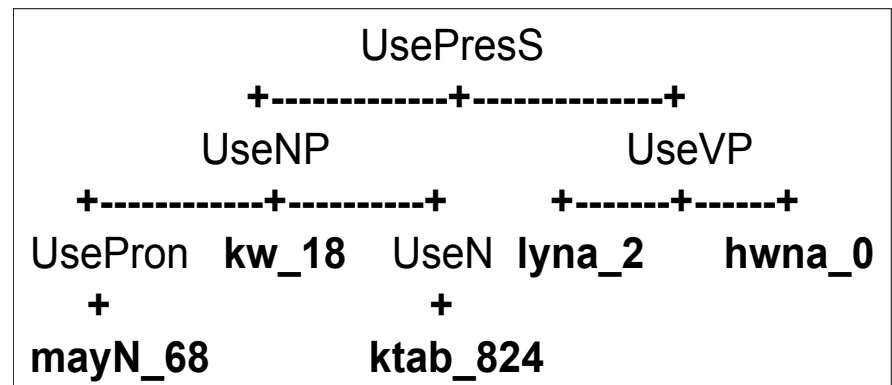
Syntactic parsing:

UsePresS

(UseNP (UsePron mayN\_68) kw\_18

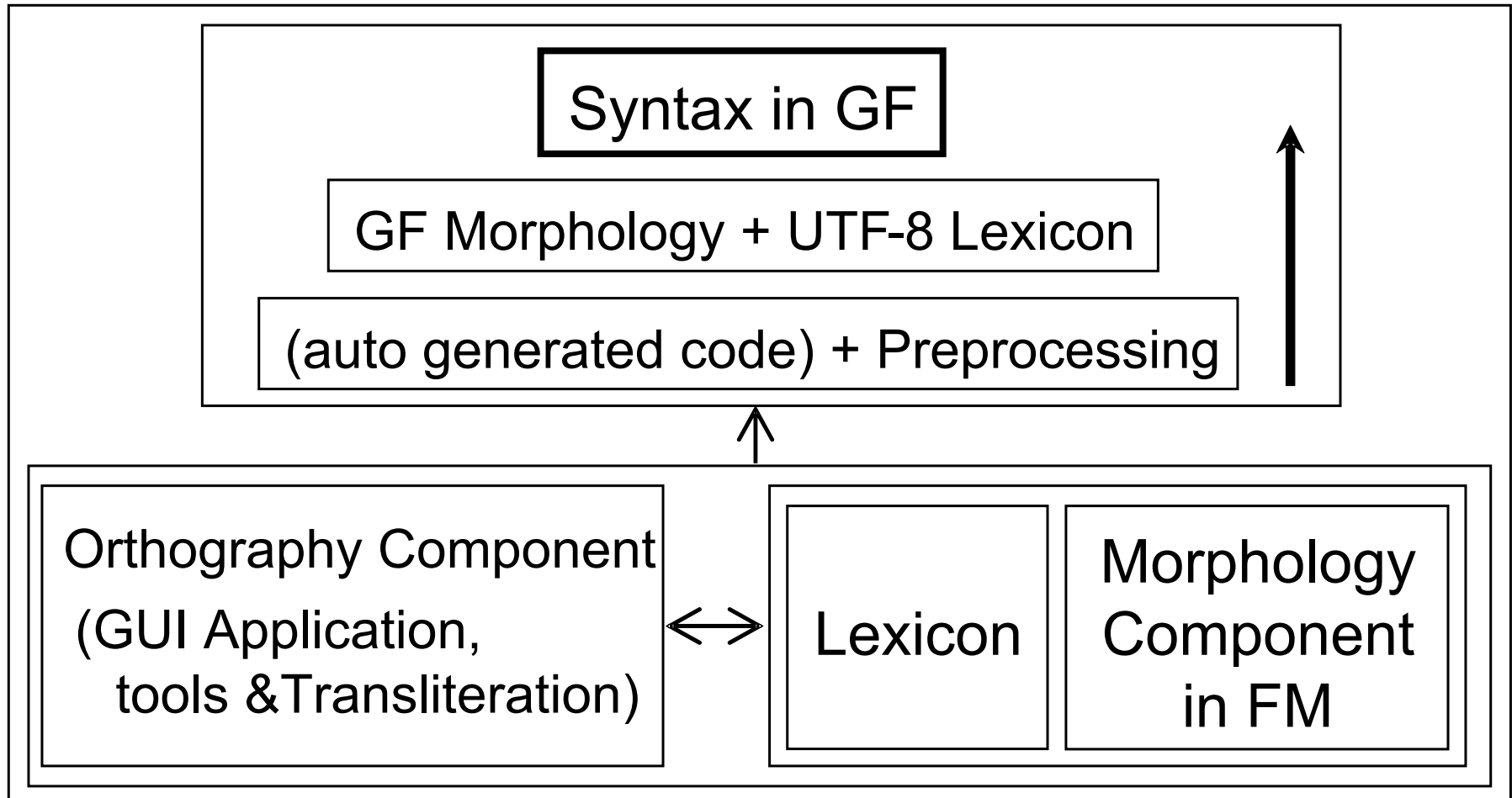
(UseN ktab\_824))

(UseVP lyna\_2 hwana\_0)



Syntax tree

# The Overall Picture



# Conclusion

## ■ Merits:

- FM proved to be a very good choice for implementing Urdu morphology
- A comprehensive, reusable & elegant implementation for Urdu that covers the linguistic abstraction (morphology) adequately

## ■ Limitations:

- A partly vocalized Lexicon
- Run time system requires an exact match:
  - One cannot check if there exist orthographically different versions of a word

# Future Work

- **A component** that matches the partly vocalized input words with the canonical words in the lexicon
- **Algorithms** to add missing diacritics on partly vocalized words
- A bigger lexicon
- A comprehensive implementation for **syntax**
- **Implementation of Hindi**, by adding a lexicon and a transliteration scheme

**Some screen shots**

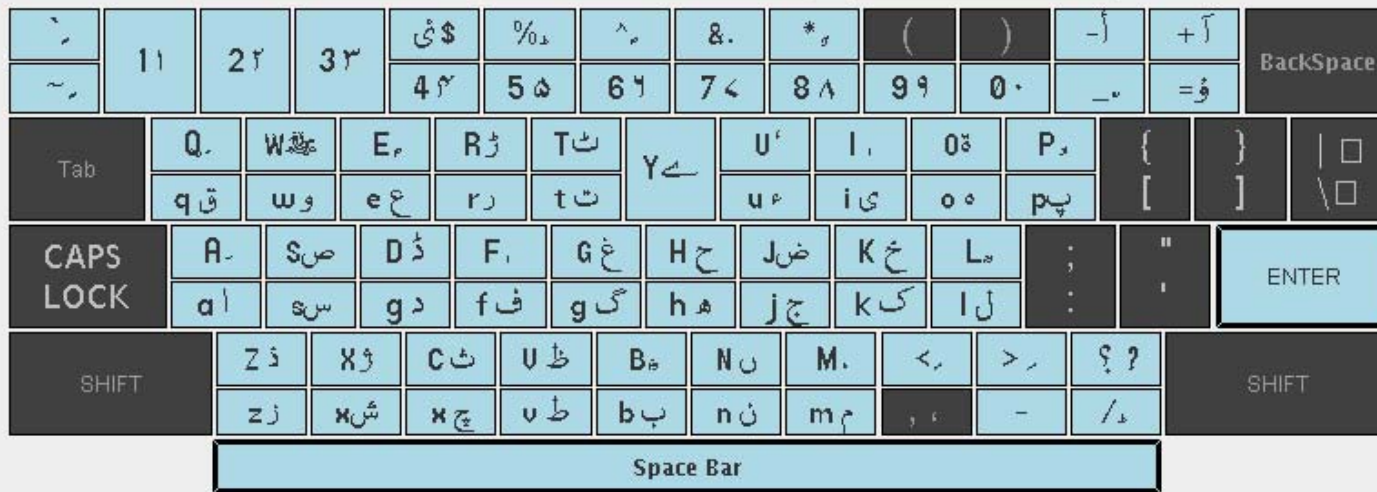
کتاب

k(i)tab

Show Analysis

Hide Urdu keyboard

Clear



Dictionary loaded: DF = 4807 and WF = 137078.

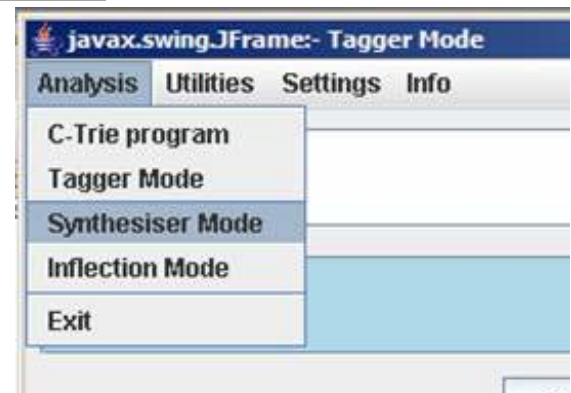
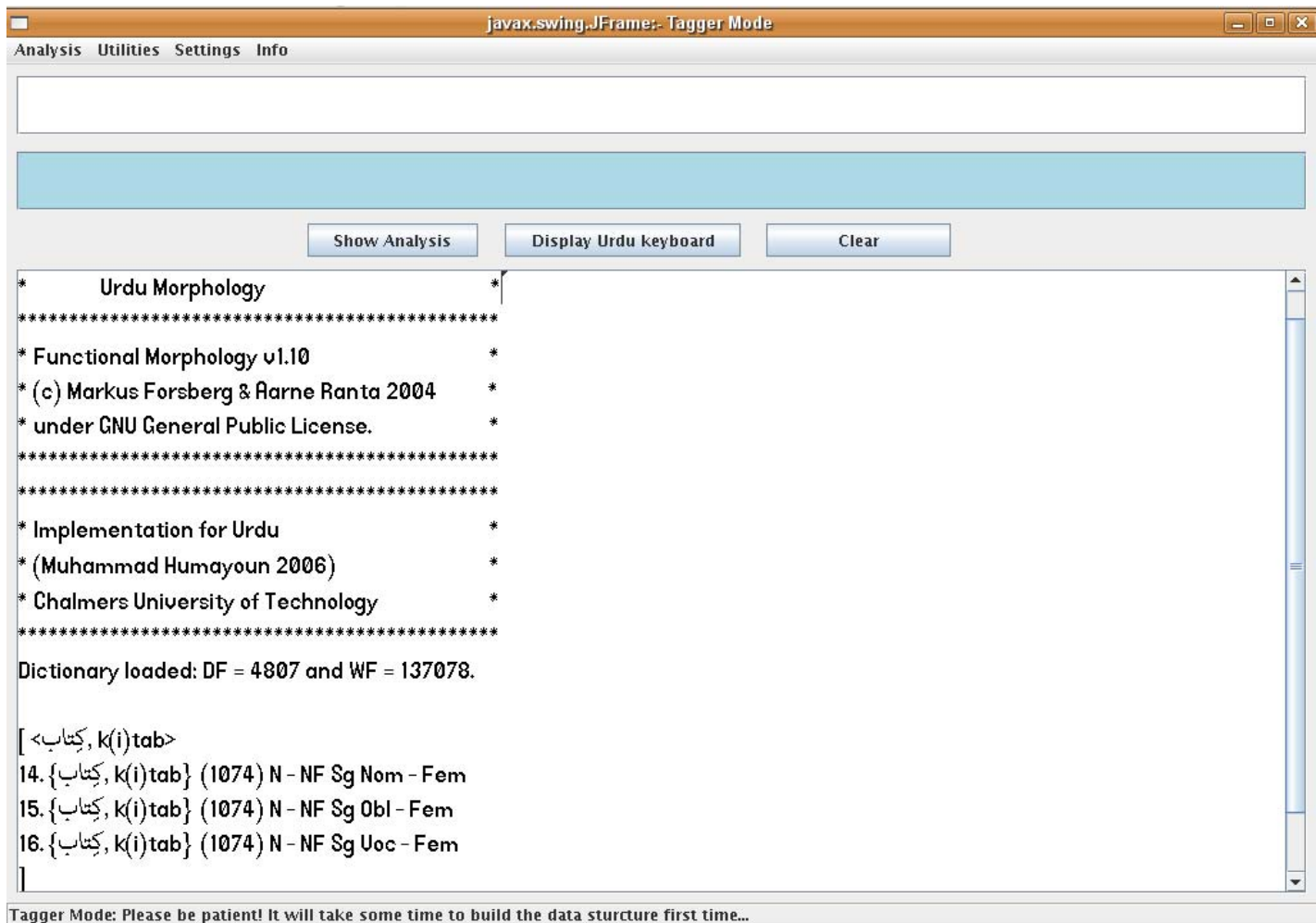
[ <کتاب, k(i)tab>

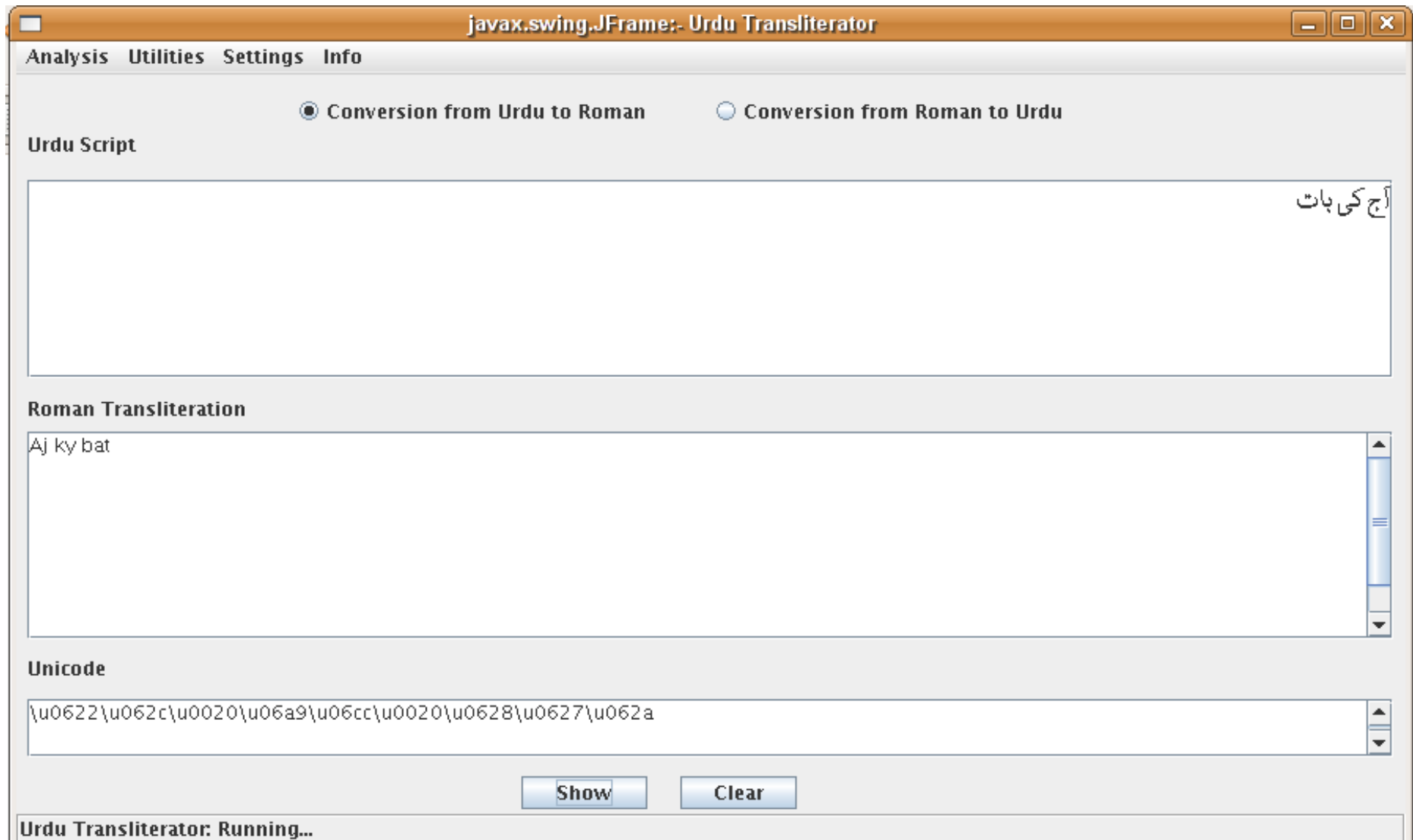
14. {کتاب, k(i)tab} (1074) N - NF Sg Nom - Fem

15. {کتاب, k(i)tab} (1074) N - NF Sg Obl - Fem

16. {کتاب, k(i)tab} (1074) N - NF Sg Voc - Fem

]









**Thanks for your attention**

Questions / Comments



Homepage of the project

<http://www.lama.univ-savoie.fr/~humayoun/UrduMorph/>

**Additional slides**

# Urdu Orthography

- Forty one non-aspirated letters:

آ اب پ ت ٹ ث ج چ ح خ د ڈ ذ ر ژ س ش ص ض ط ظ ع غ ف ق ک گ ل م ن ں و ہ ء  
ی ے

- Fifteen aspirated letters:

بھ پھ تھ ٹھ جھ چھ دھ ڈھ رھ ژھ کھ گھ مھ لھ نہ

- Three hamzah carrier (The glottal stop): اؤئی

- The Vowels & other diacritics (Aerab / Harkat, اعراب/حرکات):

ہ	a
ہ	i
ہ	u

ہ	a
ہ	i
ہ	u

ہ	ən
ہ	in
ہ	un

# Urdu Orthography - Transliteration

- Transliteration is a strict, reversible, one to one string mapping from one system of writing into another.
- Each Unicode value of Urdu alphabet is mapped with a unique Roman string
- An attempt to make transliteration as phonetic as possible
- An open source API (ICU4J) is used which is developed by IBM
- It provides a *Transliterator class* for this purpose

# Urdu Orthography - Transliteration

```
public class UrduUnicode
{
    public static final char alif='\u0627';
    public static final char bay='\u0628';
    public static final char pay='\u067e';
    .....
}
```

```
public class UrduRoman
{
    public static final String alif = "a";
    public static final String bay = "b";
    public static final String pay = "p";
    ....
}
```

```
private static final String unicode_to_Roman_rules =
    UrduUnicode.alif + ">" + UrduRoman.alif + ";" +
    UrduUnicode.bay + ">" + UrduRoman.bay + ";" +
    .....
```

```
public static Transliterator roman_to_unicode =
    Transliterator.createFromRules("RomanUrdu-Unicode", roman_to_Unicode_rules, 0);
```

```
String romanText =
    Transliterator_ur.unicode_to_roman.transliterate("Unicode Text");
```

# Urdu Orthography - Transliteration

Examples:

کتاب	Kiṭāb	k(i)tab	book
کوشش	koshish	k(a)wX(i)X	struggle
بلاؤ	bulaʔu	b(o)law^	to call

ک	→	k	اَ	→	(a)
کِ	→	(i)	ش	→	X
ت	→	t	اُ	→	(o)
ا	→	a	ؤ	→	w^
ب	→	b			

# Other word classes

## ■ Adjectives

### ■ Marked (Inflects in number, case and gender)

■ Ends with ( ا , a): nīla نیلا, nīlī نیلی, nīle نیلے : Blue

■ Ends with ( ان , an): ḡayaṅ دایان, ḡaʔeṅ دائیں, ḡaʔiṅ دائیں : Right

### ■ Unmarked

■ No inflection : khush خوش : Happy

■ Inflects in degree (Persian's style)

■ bəḡ بد, bəḡtr بدتر, bəḡtrin بدترین : bad, worse, worst

## ■ Adverbs

## ■ The closed classes

■ Pronouns, PostPositions, Particles, Interjunctions, Conjunctions, Negations, Questions and Numerals



# Lexicon Extraction – Paradigms

- Singular Feminine nouns not ending with ( ا , a), ( ن , N), ( و , w)  
( كتاب , Kiṭab, book), ( گاجر , gadʒər, carrot)

regex Not\_awN = char\* (char- ("a" | "N" | "w"));

```
paradigm n9 [x:Not_awN] =
  x { (x &
      (x+"yN" | x+"wN" | x+"w"))
    };
```

	Nom	Oblique	Voc
Sg	كتاب Kiṭab		
Pl	كتابين Kiṭaben (yN)	كتابون Kiṭabon (wN)	كتابو Kiṭabo (w)