

Challenges of Processing South Asian Languages (CPSAL)

Kengatharaiyer **Sarveswaran**

University of Konstanz, Germany

University of Jaffna, Sri Lanka.

Tafseer Ahmed

Senior NLP Scientist

Alexa Translations

Urdu/Hindi Morphology

Urdu(/Hindi) Noun

Form	Number	Gender		
Nominative	Singular	Masculine	لڑکا آیا	laRk-ā āyā boy-M.Sg.Nom came
Nominative	Plural	Masculine	دو لڑکے آئے	do laRk-ē āē two boy-M-Pl.Nom came
Oblique	Singular	Masculine	لڑکے نے کہا	laRk-ē=nē kahā boy-M.Sg.Obl=ērg said
Oblique	Plural	Masculine	لڑکوں نے کہا	laRk-oN=nē kahā boy-M.Pl.Obl=Erg said

Urdu/Hindi Noun Morphology

	Transl.	Sg.Nom	Sg.Obl	Pl.Nom	Pl.Obl
Masc. ending with ā sound	<i>boy</i>	laRk-ā لڑکا	laRk-ē لڑکے	laRk-ē لڑکے	laRk-ON لڑکوں
Masc. Other endings	<i>house</i>	gHar گھر	gHar گھر	gHar گھر	gHar-ON گھروں
Fem. ending with ī sound	<i>girl</i>	laRk-ī لڑکی	laRk-ī لڑکی	laRk-īāN لڑکیاں	laRk-īON لڑکیوں
Fem. Other endings	<i>book</i>	kitāb کتاب	kitāb کتاب	kitāb-ēN کتابیں	kitāb-ON کتابوں

Borrowed Morphology (Urdu)

Singular Word	Language of Origin	Native Plural	Borrowed Plural
darj-ah درجہ	Arabic	darj-ē درجے	darj-āt درجات
kampiuTar کمپیوٹر	English	kampiūTar-ON کمپیوٹروں	kampiūTar-z کمپیوٹرز
kitāb کتاب	Arabic	kitāb-ēN کتابیں	kutub کتب

Case Markers (Urdu/Hindi)

لڑکی نے لڑکے کو دوربین سے دیکھا

laRki=nē

laRkē=ko

dūrbīN=sē

dēkHA

girl.F.Sg.Obl=Erg boy.M-Sg.Obl=Acc telescope.F.Sg.Obl=inst saw

'The girl saw the boy with the telescope.'

Case	Case Marker
Ergative	nē
Accusative, Dative	ko
Ablative, Instrument	sē
Locative	mēN ('in'), par ('on'), ...
Genitive	k- (kā, kī, kē)

Morphological Case (Sindhi)

<http://www.languagesgulper.com/eng/Sindhi.html>

Nominative	Oblique	Ablative
gHar-u house.M-Sg.Nom	gHar-a house.M-Pl.Obl	gHar-āN house.M-Sg.Abl
gHar-āN house.M-Pl.Nom	gHar-ani house.M-Pl.Obl	gHar-ani-āN house.M-Pl.Obl-Abl

ياسين گهر کان آيو

Yasīn gHar=kHāN āyo
Yaseen house=Abl came
'Yaseen came from the house.'

ياسين گهران آيو

Yasin gHar-āN āyo
Yaseen house-Abl came
'Yaseen came from the house.'

Morphological Case (Punjabi)

- gHar-ON
home.M-Abl
- gHar-iC
home.M-Loc_in

Urdu/Hindi Verb Morphology

Word	Transliteration	Features
لکھ	likH	write
لکھتا	likH-tā	write-impf.M.Sg
لکھتی	likH-tī	write-īmpf.F.Sg
لکھی	liKH-ī	write-Perf.M.Sg
لکھوں	likH-ON	write-Subj.1P.Sg
لکھیے	likH-iē	write-Perc.2P

Urdu/Hindi Verb Morphology

Imperfective	likH-tā M.Sg	likH-tē M.PI	likH-tī F.Sg	likH-tīN F.PI
Perfective	likH-ā M.Sg	likH-ē M.PI	likH-ī F.Sg	likH-īN F.PI
Subjunctive	likH-ūN 1P.Sg	likH-ē 1P.Sg/2P.Sg	likH-ēN PI	
Future	likH-ūN-gā 1P.M.Sg likH-ūN-gī 1P.F.Sg	likH-ē-gA 2/3P.M.Sg likH-ē-gī 2/3P.F.Sg	likH-ēN-gē M.PI likH-ēN-gī F.PI	
Percative	likH-iē (2P.PI)			

Urdu/Hindi Verb Morphology

laRk-ī	ā-ī
girl-F.Sg.Nom	come-Perf.F.Sg
'A/The girl came.'	

laRk-īāN	ā-īN
girl-F.Pl.Nom	come-Perf.F.Pl
'Girls came.'	

ammi	ā-īN
mother. Hon -F.Sg.Nom	come-Perf.F. Hon
'Mother came.'	(Honorific)

Pronominal Suffixes (Sindhi)

Root: puT-u (son-M.Sg)

Suffix Person	Suffix Number	Word	Alternate phrase
1	Sg	پڻم puTu-mi (my son)	منهنجو پڻ munHanjO puTu
1	Pl	پڻون puTu-ūN (our son)	اسان جو پڻ Asān=jo puTu
2	Sg	پڻين puTu-ēN (your.sg son)	توهان جو پڻ tuHān=jo puTu
2	Pl	پڻوا puTu-va (your.pl son)	توهان جو پڻ tuHān=jo puTu
3	Sg	پڻس puTu-si (his son)	هن جو پڻ hun=jo puTu
3	Pl	پڻن puTu-ni (their son)	سندن پڻ sundanu puTu

Pronominal Suffixes (Sindhi)

on Case Markers

هن مون کي ڪتاب ڏنو

Huni mūN=kHē

he.Obl 1P.Sg.Obl=Dat

He gave me a book.

kitābu

book.Sg.M

dinū

give.PastPart.M.Sg

هن کيم ڪتاب ڏنو

Huni kHē-mi

3P.M.Obl Dat-1P.Sg

He gave me a book.

kitAbu

book.Sg.M

dinU

give.PastPart.M.Sg

Pronominal Suffixes (Sindhi)

On Verbs

خط لکيو مانس

xat-u

likH-iyO-maan-si

Letter.M-Nom.Sg

write-PastPart.Sg-1P.Sg-3P.Sg

I wrote him/her (a) letter.

Pronominal Suffixes

<https://pr.hec.gov.pk/jspui/bitstream/123456789/10182/1/PhD%20Thesis%20%28Mutee%20u%20Rahman%29%20.pdf>

http://www.sindhiaadabiboard.org/Catalogue/Lasaniyat/Book5/Book_page7.html

Tokenization

Urdu script

an-ginat	ان گنت	uncountable
an-kahī	ان کہی	untold
aasaan-tareen	آسان ترین	easiest

Tokenization

Word With Space	Transliteration and gloss	Meaning
ان گنت	an-ginat not-counted	uncountable
ان کھی	an-kahl not-told	untold
انجان	an-jAn not-known	unknown
اہم ترین	aham-tarIn important-most	most important
شادی شدہ	SHAdI-SHudah marriage-became	married.Adj
گم شدہ / گمشدہ	gum-SHudah lost-became	lost.Adj

Tokenization

Stuttgart

Transliteration	Google Search Results
شٹٹگارٹ	15
شٹٹ گارٹ	39

The rule(s) for space insertion/deletion (including the exceptions) is/are yet to be explored.

Morphology Induction

- **Linguistica**

- <https://linguistica-uchicago.github.io/lxa5/>
- Unsupervised Learning from Corpus
- Learning
- Affixes, Signatures, and associated words
- An Example is in the Python notebook

Morphology Induction

- **Morfessor**

- <https://morfessor.readthedocs.io/>
- Pre-trained segmentation models
- Models can be trained

Morphological Features

UniMorph

Number	Greater paucal	GPAUC
Number	Greater plural	GRPL
Number	Inverse	INVN
Number	Paucal	PAUC
Number	Plural	PL
Number	Singular	SG
Gender	Bantu Noun Classes	BANTU1-23
Gender	Feminine	FEM
Gender	Masculine	MASC
Gender	Nakh-Daghestanian Noun Classes	NAKH1-8
Gender	Neuter	NEUT

<https://unimorph.github.io/>

Tamil morphological analyser

<https://nlp-tools.uom.lk/thamizhi-morph/parse-sentence.php>

Thank you

Tafseer Ahmed

Senior NLP Scientist, Alexa Translations

tafseer@gmail.com

Kengatharaiyer Sarveswaran

University of Jaffna, Sri Lanka.

University of Konstanz, Germany

sarves@univ.jfn.ac.lk

[sarves.github.io](https://github.com/sarves)

Subwords for Deep Learning

What are the challenges of rich morphology to:

Morphological Analysis/Generation

Part of Speech tagging

Word Embedding

Morphological Features

ParGram(Parallel Grammar)

Nominal Features

ANIM: \pm

CASE: ACC, DAT, GEN, NOM

GEND: FEM, MASC, NEUT

NTYPE: COUNT, MASS, PROPER

NUM: PL, SG

PERS: 1, 2, 3

REFL: \pm

Morphological Features

Universal Dependency

Lexical features	Inflectional features	
	<i>Nominal*</i>	<i>Verbal*</i>
<u>PronType</u>	<u>Gender</u>	<u>VerbForm</u>
<u>NumType</u>	<u>Animacy</u>	<u>Mood</u>
<u>Poss</u>	<u>NounClass</u>	<u>Tense</u>
<u>Reflex</u>	<u>Number</u>	<u>Aspect</u>
<u>Foreign</u>	<u>Case</u>	<u>Voice</u>
<u>Abbr</u>	<u>Definite</u>	<u>Evident</u>
	<u>Degree</u>	<u>Polarity</u>
		<u>Person</u>
		<u>Polite</u>
		<u>Clusivity</u>

<https://universaldependencies.org/u/fe>