

Developing Standards and Linguistic Resources for Computational Research in Pakistani Languages



www.CLE.org.pk

Sarmad Hussain
 Professor
 Center for language Engineering (CLE)
 Al-Khwarizmi Institute of Computer Sciences (KICS)
 University of Engineering and Technology (UET) Lahore
sarmad.hussain@kics.edu.pk

1

Need

ICTs promise significant socio-economic impact
 Impact dependent on size of population which can use ICTs

- 180 Million citizens need access
- 66+ languages
- 10% understand English
- 58% literate
- 11% have access to computers
- 70% have access to mobile phones
- ITU IDI: Pakistan ranked 127 of 155 nations

Human Language Technology necessary to bridge the gap

www.cle.org.pk

2

Languages of Pakistan

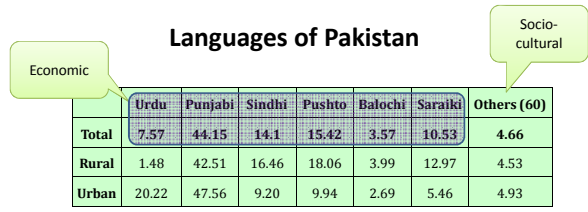
	Urdu	Punjabi	Sindhi	Pushto	Balochi	Saraiki	Others (60)
Total	7.57	44.15	14.1	15.42	3.57	10.53	4.66
Rural	1.48	42.51	16.46	18.06	3.99	12.97	4.53
Urban	20.22	47.56	9.20	9.94	2.69	5.46	4.93

Percent Population of Pakistan by Mother Tongue

www.cle.org.pk

3

Languages of Pakistan



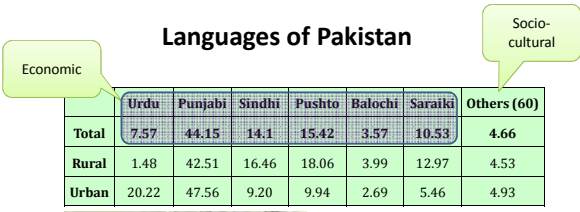
	Urdu	Punjabi	Sindhi	Pushto	Balochi	Saraiki	Others (60)
Total	7.57	44.15	14.1	15.42	3.57	10.53	4.66
Rural	1.48	42.51	16.46	18.06	3.99	12.97	4.53
Urban	20.22	47.56	9.20	9.94	2.69	5.46	4.93

Percent Population of Pakistan by Mother Tongue

www.cle.org.pk

4

Languages of Pakistan



	Urdu	Punjabi	Sindhi	Pushto	Balochi	Saraiki	Others (60)
Total	7.57	44.15	14.1	15.42	3.57	10.53	4.66
Rural	1.48	42.51	16.46	18.06	3.99	12.97	4.53
Urban	20.22	47.56	9.20	9.94	2.69	5.46	4.93

Percent Population of Pakistan by Mother Tongue



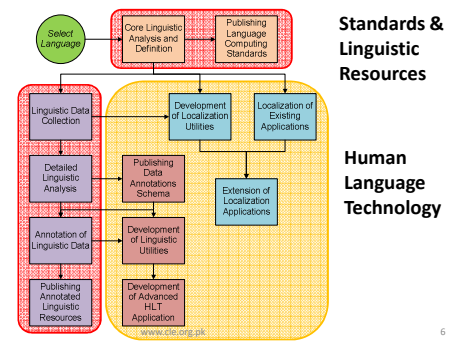
Languages of Pakistan in Danger (UNESCO) (23)

- Vulnerable
- definitely endangered
- severely endangered

www.cle.org.pk

5

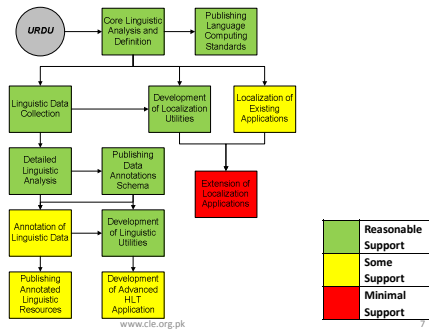
Development Process of Human Language Technology



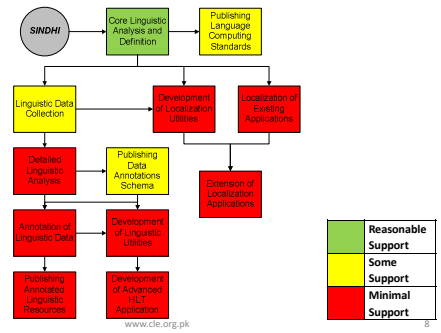
www.cle.org.pk

6

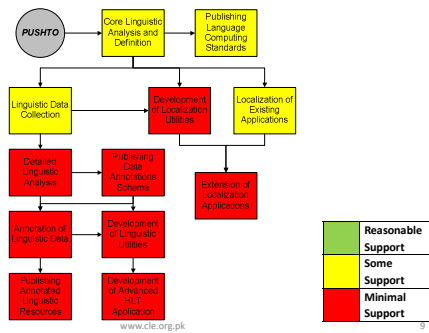
Status of Human Language Technology



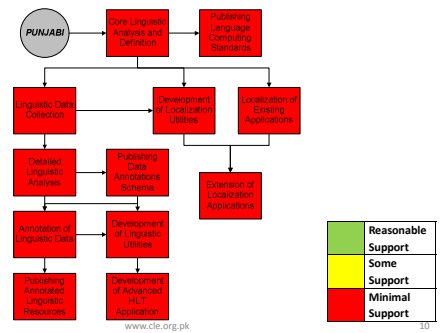
Status of Human Language Technology



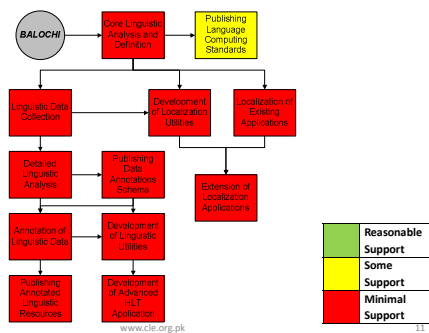
Status of Human Language Technology



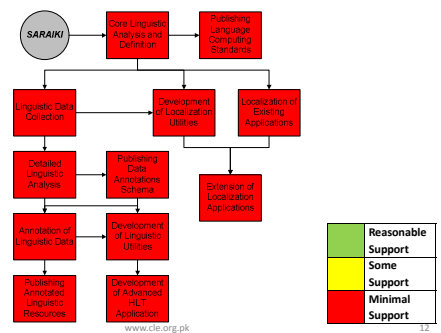
Status of Human Language Technology



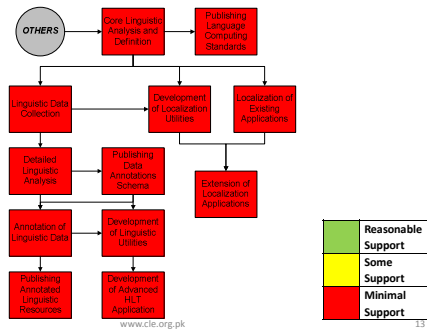
Status of Human Language Technology



Status of Human Language Technology



Status of Human Language Technology



Core Linguistic Analysis and Definition

- Linguistic details incomplete or unavailable
- Relevant cultural conventions rarely documented
- Need analysis of language and cultural conventions
 - Should be driven by linguists and community
- Precise definitions of relevant aspects
 - Involvement of technologists to iron out ambiguities

Core Linguistic Analysis and Definition

- Character-set
 - Alphabet – letters & aerab
 - Punctuation marks
 - Special symbols, etc.
- Sorting or collation of characters
- Cultural conventions for representing
 - Numbers
 - Time
 - Calendar
 - Currency
- Translation of terms for software interface

Standards

- Character Encoding – ISO 10646 - Unicode
- Locale – ISO 639 – ISO 3166 – Unicode-CLDR
 - Keyboard
 - Collation Element
- Localization Terminology
- Develop national standards
- Follow up with international standards

Lexica

- Word lists
- Annotated Lexical
 - Pronunciation
 - Agreement information (num, gender, case, respect)
 - Part of Speech (POS)
 - Translation
 - Sub-categorization frame/argument structure
 - Root
 - Morphological information
- Wordnet

Word Lists

- Word list – surface forms
- Lexeme list
- Closed class words
- Proper Names
- Frequency based word lists

Ligature	Frequency	Word	Frequency
ا	2885748	کے	743842
ر	1716383	میں	852882
و	1350586	کو	515545
د	896166	ہیں	466800
کے	746801	ہو	415700
میں	582882	تو	368355
کی	581999	ہے	306300
سے	508375	کی	281862
ن	452167	میں	254365
کے	448244	ہیں	244817
کا	447729	کو	237419
تو	424334	ہیں	217840
سے	388097	ہے	207801
کی	362535	کو	173467
ن	348302	میں	148888
کی	342672	ہو	140387
		تو	128900
		ہے	115241
		میں	113465
		کو	88837
		ہیں	87615
		تو	86765
		ہے	84582

Phonetic Lexicon

Attock	A T T A K
Bahawalpur	B A H A A V A L P U R
Bhakar	B _ H A K A R
Chakwal	T _ S H A K V A A L
Faisalabad	F A Y S L A B A A D _ D
Gujranwala	G U D _ Z Z R A A N V A A L A A N
Gujrat	G U D _ Z Z R A A T T
Jhang	D _ Z Z _ H A N G
Jhelum	D _ Z Z A E H L A M
Kasur	K A S U U R
Khushab	X U U S H A A B
Lahore	L A A H O R
Mianwali	M I I A A N V A A L I
Multan	M U L T _ D A A N
Rawalpindi	R A A V A L P I N D D I I
Sahiwal	S A A H I I V A A L
Sargodha	S A R G O O D _ D _ H A A
Sheikhpura	S H A E X U U P U R A A
Sialkot	S I A L K O O T T

www.cle.org.pk 19

Urdu ID	English ID	English Word	Category	Concept	Example	Synset
100853	6854161	English	Noun	انگلینڈ کی زبان	اسے عربی اور فارسی کے ساتھ ساتھ انگریزوں کی پہلی، انگلستان کی انگریزی، انگلش	
103196	617535	reading	Verb	مطالعہ کرنا	آج یہ کتاب پڑھنے کا موقع ملا ہے	دیکھنا، پڑھنا، مطالعہ کرنا
101288	4977386	intensity	Noun	شہید ہونے کی خصوصیت	ہوا اتنی زور سے ملی کہ خیر اڑ گیا	بھوش، زور، شدت، ہمتی
103802	2179176	give	Verb	مہینہ بیسہ امتت رکھنا	اس نے اپنے روپے اس کو سوچے	دینا، سونپنا
102950	145831	certainly	Adverb	بے شک کسی شے کے بغیر	یقیناً اس کا کام منکل ہو گیا ہو گا	یاقین، بلاشبہ، بے شک، ضرور، یقیناً
104727	9771320	child	Noun	اسٹیشن پر دم دھاکے کا زخمی سات	سارے بچے اسپتال میں دم توڑ گیا	بچہ، طفل

www.cle.org.pk 20

Text Corpora

- Text corpus
 - Size
 - Time
 - Genre
- CLE Urdu Digest corpus
 - Size
 - 100k
 - 500k
 - 1M
 - Time
 - 2003-2012
 - Genres
 - Balanced

Category	Sub-category	Percentages
1. Informational (80%)		
a) Informal	Letters	10%
	Interviews	10%
a) Formal		
	Press	8%
	Religion	8%
	Sports	8%
	Culture (travel, history)	8%
	Entertainment	4%
	Health	8%
	Science (education, technology)	16%
1. Imaginative (20%)		
	Short Stories	8%
	Translation of foreign literature	4%
	Novels	4%
	Book reviews	4%

www.cle.org.pk 21

Text Corpora

- State of the Art
 - English: Web 1T 5-gram Version 1
 - Published by Google
 - Number of tokens: **1,024,908,267,229**
 - Number of sentences: 95,119,665,584
 - Number of unigrams: 13,588,391
 - Number of bigrams: 314,843,401
 - Number of trigrams: 977,069,902
 - Number of fourgrams: 1,313,818,354
 - Number of fivegrams: 1,176,470,663

www.cle.org.pk 22

English: Web 1T 5-gram Version 1

- Sample Trigrams and Frequencies
 - ceramics community for 61
 - ceramics companies . 53
 - ceramics companies consultants 173
 - ceramics company ! 4432
 - ceramics company , 133
 - ceramics company . 92
 - ceramics company 41
 - ceramics company facing 145
 - ceramics company in 181
 - ceramics company started 137
 - ceramics company that 87
 - ceramics component (76
 - ceramics composed of 85
 - ceramics composition as 41
 - ceramics computer graphics 51
 - ceramics computer imaging 52
 - ceramics consist of 92

www.cle.org.pk 23

Annotated Corpora

- Text Corpus
 - POS tagged
 - Chunk tagged
 - Sense tagged
- Tree Bank
 - Syntax
 - Grammatical relations
 - Dependency relations

www.cle.org.pk 24

دنیا کا ہر فرد کامیابی کا آرزومند ہے۔ ناکامی سے سب گھبراتے ہیں۔ عزت، دولت، راحت اور عافیت کی زندگی کے سبھی شیدائی ہیں۔ لیکن اصل کامیابی کیا چیز ہے؟ اور حقیقی عزت و راحت کس طرح نصیب ہوتی ہے؟ اس مجید سے بہت کم لوگ واقف ہیں۔ اگر آپ حقیقی کامیابی کے گرہ پائنا چاہتے ہیں تو ڈاکٹر زاہد منیر مہار کی تازہ تصنیف آئیڈیہ کروڈا پڑھیے۔ ۱۱۳ صفحات کی اس کتاب کا ایک ایک حرف بصیرت کے در سے کھینچنے پر مامور ہے۔

دنیار / ہر / جلد / کامیابی / ہر / آرزومند / ہے /
 ہر / کامی / ہے / جلد / گھبراتے / ہیں /
 راحت / اور / عافیت / کی / زندگی / کے / سبھی / شیدائی / ہیں /
 لیکن / اصل / کامیابی / کیا / چیز / ہے /
 اور / حقیقی / عزت / و / راحت / کس / طرح / نصیب / ہوتی / ہے /
 اس / مجید / سے / بہت / کم / لوگ / واقف / ہیں /
 اگر / آپ / حقیقی / کامیابی / کے / گرہ / پائنا / چاہتے / ہیں /
 تو / ڈاکٹر / زاہد / منیر / مہار / کی / تازہ / تصنیف / آئیڈیہ / کروڈا / پڑھیے /
 ۱۱۳ / صفحات / کی / اس / کتاب / کا / ایک / ایک / حرف / بصیرت / کے / در / سے / کھینچنے / پر / مامور / ہے /

www.cle.org.pk

```
( (S
(NP Battle-tested industrial managers
here)
always
(VP buck
up
(NP nervous newcomers)
(PP with
(NP the tale
(PP of
(NP (NP the
(ADJP first
(PP of
(NP their countrymen)))
(S (NP *)
to
(VP visit
(NP Mexico))))
(NP (NP a boatload
(PP of
(NP (NP warriors)
(VP-1 blown
ashore
(ADVP (NP 375 years
ago))))
(VP-1 *pseudo-attach*)))))))))
```

- ((S (NP Battle-tested industrial managers here) always (VP buck up (NP nervous newcomers) (PP with (NP the tale (PP of (NP (NP the (ADJP first (PP of (NP their countrymen))) (S (NP *) to (VP visit (NP Mexico))))), (NP (NP a boatload (PP of (NP (NP warriors) (VP-1 blown ashore (ADVP (NP 375 years ago)))))) (VP-1 *pseudo-attach*)))))) .)

Sentence Count: 317 Displayed Tree (Sentence): 9

Arabic (Quran)

Chapter (97) surah Al-Mulk (Dominant)

Phrase Structure vs. Dependency Treebanks

www.cle.org.pk

Sense Tagged Corpus

Text Corpora

- Parallel Text corpus
 - Language
 - Size
 - Genre
 - Time
 - Alignment
- Spelling Error corpus

Parallel Corpus

- English

Pierre Vinken , 61 years old , will join the board as a nonexecutive director Nov. 29 .

Mr. Vinken is chairman of Elsevier N.V. , the Dutch publishing group .
- Urdu

آکٹو ۲۹ نوں ۶۱ سالوں وچ پیر وینکن ، 61 سالوں وچ ، ایگزیکٹو ڈائریکٹر کے طور پر بورڈ میں شامل ہونے کے لیے ایک غیر ایگزیکٹو ڈائریکٹر کے طور پر نوں ۲۹ نوں .

مسٹر وینکن ایگزیکٹو ڈائریکٹر کے طور پر الٹیر ایف ڈی ، ڈچ پبلشنگ گروپ کے چیئرمین ہیں .
- Nepali

६१ <CD> वर्षीय <I> पियरे <NNP> विन्केन <NNP> नोभेम्बर <NNP> २९ <CD> बाट <POP> सल्लाहकार <NN> को <PKO> रूप <NN> मा <POP> सञ्चालक <NN> समिति <NN> मा <POP> आउनुहुनेछ <VBX> । <YF>

श्री <NN> विन्केन <NNP> डच <NNP> प्रकाशन <NN> समूह <NN> एल्सेवियर <NNP> एन. वी. <FB> को <PKO> अध्यक्ष <NN> हुनुहुन्छ <VBF> । <YF>

Corpora

- Speech Corpus
 - Duration
 - Speakers
 - Channels
 - Environments
 - Contents
 - Accent(s)
 - Gender
 - Style (read/spontaneous)
 - Vocabulary

Annotated Corpora

- Speech Corpus
 - Textual transcription
 - Segmental IPA transcription
 - Aligned transcription
 - Syllable boundaries
 - Stress
 - Intonation
 - POS

Speech Tagging Schema

p1^p2-p3+p4=p5@p6 p7
 /A:a1 a2 a3
 /B:b1-b2-b3@b4-b5&b6-b7#b8-b9\$b10-b11|b12-b13;b14-b15|b16
 /C:c1+c2+c3
 /D:d1 d2 /E:e1+e2@e3+e4&e5+e6#e7+e8 /F: f1 f2
 /G:g1 g2 /H:h1=h2^h3=h4|h5 /I:i1 i2

p1	the phoneme identity before the previous phoneme
p2	the previous phoneme identity
p3	the current phoneme identity
p4	the next phoneme identity
p5	the phoneme after the next phoneme identity
p6	position of the current phoneme identity in the current syllable (forward)
p7	position of the current phoneme identity in the current syllable (backward)

s1	whether the previous syllable stressed or not (0: not stressed, 1: stressed)
s2	whether the previous syllable accented or not (0: not accented, 1: accented)
s3	the number of phonemes in the previous syllable
s4	whether the current syllable stressed or not (0: not stressed, 1: stressed)
s5	whether the current syllable accented or not (0: not accented, 1: accented)
s6	the number of phonemes in the current syllable
s7	position of the current syllable in the current word (forward)
s8	position of the current syllable in the current word (backward)
s9	position of the current syllable in the current phrase (forward)
s10	position of the current syllable in the current phrase (backward)
s11	the number of stressed syllables before the current syllable in the current phrase
s12	the number of accented syllables before the current syllable in the current phrase
s13	the number of accented syllables after the current syllable in the current phrase
s14	the number of stressed syllables after the current syllable in the current phrase
s15	the number of syllables from the previous stressed syllable to the current syllable
s16	the number of syllables from the previous accented syllable to the current syllable
s17	the number of syllables from the current syllable to the next accented syllable
s18	the number of syllables from the current syllable to the next stressed syllable
s19	name of the vowel of the current syllable
t1	whether the next syllable stressed or not (0: not stressed, 1: stressed)
t2	whether the next syllable accented or not (0: not accented, 1: accented)
t3	the number of phonemes in the next syllable

a ₁	gpos (guess part-of-speech) of the previous word
a ₂	the number of syllables in the previous word
a ₃	gpos (guess part-of-speech) of the current word
a ₄	the number of syllables in the current word
a ₅	position of the current word in the current phrase (forward)
a ₆	position of the current word in the current phrase (backward)
a ₇	the number of content words before the current word in the current phrase
a ₈	the number of content words after the current word in the current phrase
a ₉	the number of words from the previous content word to the current word
a ₁₀	the number of words from the current word to the next content word
f ₁	gpos (guess part-of-speech) of the next word
f ₂	the number of syllables in the next word
b ₁	the number of syllables in the previous phrase
b ₂	the number of words in the previous phrase
k ₁	the number of syllables in the current phrase
k ₂	the number of words in the current phrase
k ₃	position of the current phrase in this utterance (forward)
k ₄	position of the current phrase in this utterance (backward)
k ₅	TOBI code of the current phrase
i ₁	the number of syllables in the next phrase
i ₂	the number of words in the next phrase

www.cle.org.pk 37

Speech Annotated File Sample

p1^p2-p3+p4=p5@p6 p7
 /A:a1 a2 a3
 /B:b1-b2-b3@b4-b5&b6-b7#b8-b9\$b10-b11|b12-b13;b14-b15|b16
 /C:c1+c2+c3
 /D:d1 d2 /E:e1+e2@e3+e4&e5+e6#e7+e8 /F: f1 f2
 /G:g1 g2 /H:h1=h2^h3=h4|h5 /I:i1 i2

1807500 3432500 pau^pau^ao+th+er@1_2/A:0_0_0/B:1-1-2@1-2&1-7#1-4\$1-310-2;0-4|ao/C:0+0+1/D:0_0/E:content+2@1+5&1+2#0+3/F:in_1/G:0_0/H:7=5^1=2|L-L%/I:7=3/I:14+8-2

3432500 4557500 pau^ao+th+er+ah@2_1/A:0_0_0/B:1-1-2@1-2&1-7#1-4\$1-310-2;0-4|ao/C:0+0+1/D:0_0/E:content+2@1+5&1+2#0+3/F:in_1/G:0_0/H:7=5^1=2|L-L%/I:7=3/I:14+8-2

www.cle.org.pk

38

Corpora

- Document Image Corpus
 - Resolution (DPI)
 - Color/Grayscale/BW
 - Print quality
 - Paper quality
 - Paper transparency
 - Genre
 - Publisher

www.cle.org.pk

39

Annotated Corpora

- Document Image Corpus
 - Textual transcription
 - Line
 - Text Area
 - Figures and page frames
 - Heading
 - Header and Footer
 - Main bodies
 - Diacritics and their association

www.cle.org.pk

40

113

ایمان، اُمید اور محبت

جبکہ عربی اور کسی حد تک اردو زبان بھی وہ بول لیتا تھا اگرچہ وہ ان زبانوں میں لکھ یا پڑھ نہیں سکتا تھا۔ اس کی اس خصوصیت کے انکشاف نے یکدم ہی اسے اپنی کلاس اور کسی حد تک اسکول میں پاپولر کر دیا تھا۔ لیکن بیچ کی کلاس میں ایک دن اتفاقاً اس کے نیچر کو اس بات کا پتا چلا تھا کہ وہ جرمن زبان پر بھی دسترس رکھتا ہے۔

”تو ڈیٹیل تم دو زبانوں کو استعمال کر سکتے ہو؟“ نیچر نے اسے سراہتے ہوئے کہا۔
 ”دو نہیں چار..... عربی اور اردو بھی۔ اگرچہ میں انہیں لکھ پڑھ نہیں سکتا مگر اس میں گفتگو کر سکتا ہوں۔“ مدہم آواز میں کہے گئے جملے نے یک دم ہی پوری کلاس کو سرموڑ کر اس کی طرف متوجہ ہونے پر مجبور کر دیا۔ ان کی آنکھوں میں حیرت کے ساتھ ساتھ سائنس بھی تھی۔

”چار زبانیں..... زبردست۔ مگر چار زبانیں کیسے؟ میرا مطلب ہے عربی اور اردو؟“
 ”میرے ڈیڈی بہت عرصے سے ڈل ایٹ اور ایشیا کے ممالک میں کام کرتے رہے ہیں، میری پیدائش بھی مراکش میں ہوئی اس لیے عربی بولنا آگئی اور پچھلے دو سال سے ہم لوگ انڈیا میں تھے۔ وہاں لوگوں سے بات چیت انگلش یا اردو میں ہی ہوتی تھی، اس لیے اس کو بھی استعمال کرنا آ گیا۔“

113

ایمان، اُمید اور محبت

جبکہ عربی اور کسی حد تک اردو زبان بھی وہ بول لیتا تھا اگرچہ وہ ان زبانوں میں لکھ یا پڑھ نہیں سکتا تھا۔ اس کی اس خصوصیت کے انکشاف نے یکدم ہی اسے اپنی کلاس اور کسی حد تک اسکول میں پاپولر کر دیا تھا۔ لیکن بیچ کی کلاس میں ایک دن اتفاقاً اس کے نیچر کو اس بات کا پتا چلا تھا کہ وہ جرمن زبان پر بھی دسترس رکھتا ہے۔

”تو ڈیٹیل تم دو زبانوں کو استعمال کر سکتے ہو؟“ نیچر نے اسے سراہتے ہوئے کہا۔
 ”دو نہیں چار..... عربی اور اردو بھی۔ اگرچہ میں انہیں لکھ پڑھ نہیں سکتا مگر اس میں گفتگو کر سکتا ہوں۔“ مدہم آواز میں کہے گئے جملے نے یک دم ہی پوری کلاس کو سرموڑ کر اس کی طرف متوجہ ہونے پر مجبور کر دیا۔ ان کی آنکھوں میں حیرت کے ساتھ ساتھ سائنس بھی تھی۔

”چار زبانیں..... زبردست۔ مگر چار زبانیں کیسے؟ میرا مطلب ہے عربی اور اردو؟“
 ”میرے ڈیڈی بہت عرصے سے ڈل ایٹ اور ایشیا کے ممالک میں کام کرتے رہے ہیں، میری پیدائش بھی مراکش میں ہوئی اس لیے عربی بولنا آگئی اور پچھلے دو سال سے ہم لوگ انڈیا میں تھے۔ وہاں لوگوں سے بات چیت انگلش یا اردو میں ہی ہوتی تھی، اس لیے اس کو بھی استعمال کرنا آ گیا۔“

ایمان، امید اور محبت

113

تک کہ عربی اور کسی حد تک اردو زبان بھی وہ بول لیتا تھا اگرچہ وہ ان زبانوں میں لکھنا نہیں سکتا تھا۔
اس کی اس خصوصیت کے انکشاف نے یکدم ہی اسے اپنی کلاس اور کسی حد تک اسکول میں پاپولر کر
دیا تھا۔ لیکن عربی کی کلاس میں ایک دن اتفاقاً اس کے بچے کو اس بات کا پتا چلا تھا کہ وہ جس زبان پر بھی
دسترس رکھتا ہے۔

”نوڈیشنل نم دوربانوں کو استعمال کر سکتے ہو؟“ بچے نے اسے سرائے ہوئے کہا۔
”دو نہیں چار..... عربی اور اردو بھی۔ اگرچہ میں انہیں لکھ پڑھ نہیں سکتا مگر اس میں گفتگو کر سکتا
ہوں۔“ مذہم آواز میں کہے گئے جملے نے یک دم ہی پوری کلاس کو سرسوز کر اس کی طرف منوجہ ہونے پر محمود
کر دیا۔ ان کی آنکھوں میں حیرت کے ساتھ ساتھ سانس بھی تھی۔

”چار زبانیں..... زبردست۔ مگر چار زبانیں کسے؟ میرا منظر ہے عربی اور اردو؟“
”میرے ڈیڑی بہت عرصے سے نڈل ایسٹ اور ایشیا کے ممالک میں کام کرتے رہے ہیں، میری
بدانت بھی مراکش میں ہوئی اس لیے عربی بولنا آگئی اور پچھلے دو سال سے ہم لوگ انڈیا میں تھے۔ وہاں
لوگوں سے بات چیت انگلش یا اردو میں ہی ہوتی تھی، اس لیے اس کو بھی استعمال کرنا آ گیا۔“




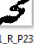














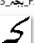
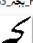
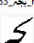
ایمان، امید اور محبت

113

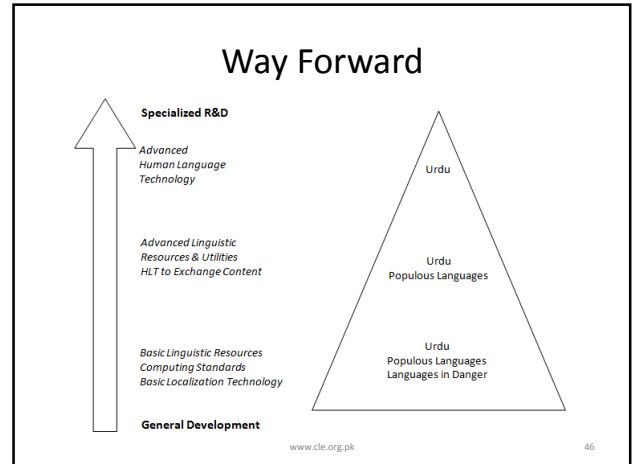
تک کہ عربی اور کسی حد تک اردو زبان بھی وہ بول لیتا تھا اگرچہ وہ ان زبانوں میں لکھنا نہیں سکتا تھا۔
اس کی اس خصوصیت کے انکشاف نے یکدم ہی اسے اپنی کلاس اور کسی حد تک اسکول میں پاپولر کر
دیا تھا۔ لیکن عربی کی کلاس میں ایک دن اتفاقاً اس کے بچے کو اس بات کا پتا چلا تھا کہ وہ جس زبان پر بھی
دسترس رکھتا ہے۔


”نوڈیشنل نم دوربانوں کو استعمال کر سکتے ہو؟“ بچے نے اسے سرائے ہوئے کہا۔
”دو نہیں چار..... عربی اور اردو بھی۔ اگرچہ میں انہیں لکھ پڑھ نہیں سکتا مگر اس میں گفتگو کر سکتا
ہوں۔“ مذہم آواز میں کہے گئے جملے نے یک دم ہی پوری کلاس کو سرسوز کر اس کی طرف منوجہ ہونے پر محمود
کر دیا۔ ان کی آنکھوں میں حیرت کے ساتھ ساتھ سانس بھی تھی۔

”چار زبانیں..... زبردست۔ مگر چار زبانیں کسے؟ میرا منظر ہے عربی اور اردو؟“
”میرے ڈیڑی بہت عرصے سے نڈل ایسٹ اور ایشیا کے ممالک میں کام کرتے رہے ہیں، میری
بدانت بھی مراکش میں ہوئی اس لیے عربی بولنا آگئی اور پچھلے دو سال سے ہم لوگ انڈیا میں تھے۔ وہاں
لوگوں سے بات چیت انگلش یا اردو میں ہی ہوتی تھی، اس لیے اس کو بھی استعمال کرنا آ گیا۔“

		
G_E_C_B10_R_P35_F14_637	G_E_C_B20_R_P95_F14_260	G_E_C_B26_R_P19_F14_1070
		
G_E_C_B51_R_P233_F14_115	G_E_C_B58_R_P245_F14_851	G_E_C_B58_R_P308_F14_368
		
G_E_C_B72_R_P166_F14_2586	G_E_C_B91_R_P11_F14_596	G_E_C_B94_R_P14_F14_491
		
G_E_C_B96_R_P8_F14_1377	G_E_C_B98_R_P3_F14_810	G_E_C_B117_R_P55_F14_332
		
G_HFL_383_بچہ_F14_CC_11	G_HFL_383_بچہ_F14_CC_14	G_HFL_383_بچہ_F14_CC_16
		
G_HFL_383_بچہ_F14_CC_22	G_HFL_383_بچہ_F14_CC_25	G_HFL_383_بچہ_F14_CC_29
		
G_HFL_383_بچہ_F14_CC_35	G_HFL_383_بچہ_F14_CC_38	G_HFL_383_بچہ_F14_CC_41

www.cle.org.pk 45




www.cle.org.pk

Thank you

Questions?

sarmad.hussain@kics.edu.pk

www.cle.org.pk 47