

Creating a Handwriting Recognition Corpus for Bushman Languages



Kyle Williams and Hussein Suleman

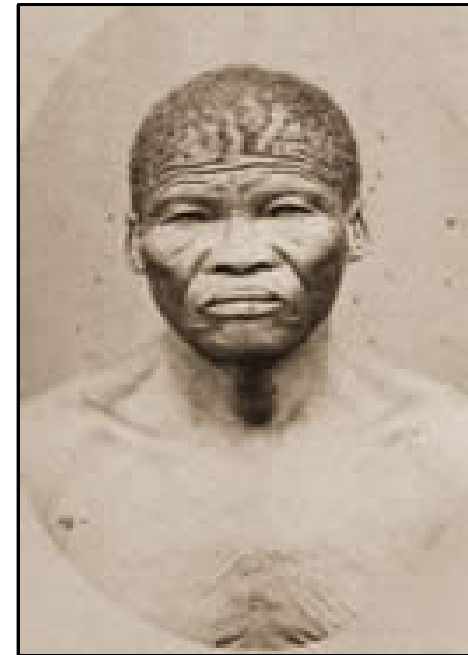




BUSHMAN PEOPLE



- Bushman people of Southern Africa
 - Earliest inhabitants of Earth
 - Unique view of the world
 - No living speakers of many Bushman languages





BLEEK AND LLOYD COLLECTION



- Collection contains notebooks, art and dictionaries
 - Bushman culture encoded in metaphorical stories
 - Preserving this collection → preserving Bushman culture



Digital Libraries Laboratory, University of Cape Town



BLEEK AND LLOYD COLLECTION



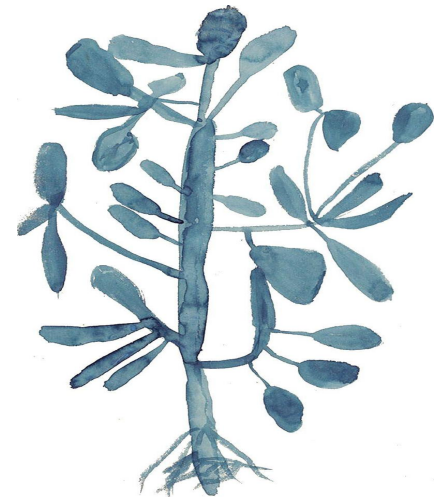
815 224
my brother, ^{tos'}abba ^a, ⁿⁱ#ha
stole (from) me ^Tne ⁺ni ^{ya}he
my quiver, ⁿⁱha [#]khewu ^a,
Nani !'ro ^a ^Tne
he ^{ka} ^{na} ⁱ, ^{na}
is foolish; he does not ^Tjebbi; ^{na} ^{pa}ki
understand, I [#]ruakhha, ⁿⁱ
understand, ⁺ka [#]ruakhha,
he strongly ^{nan} !'karsken
takes away ^{karop}, ^{butten} ⁺ ^{hin} !'nin,
he says to the other ^{nan} ⁺ ^{ka} ^{ka} !'ru
man, 'Bring me ^{ro}, ^Tne ^Tni ^{sot}
he



BLEEK AND LLOYD COLLECTION



- Already have systems for preservation and viewing collection
- Next step involves enhancing use
 - Make text searchable
 - Index text
 - Reprint of text in books
 - Text-to-speech
- Need a corpus of transcriptions





BUSHMAN TEXT



- Text contains complex diacritics
 - Stacked above and below characters
 - Span multiple characters

!hu i' ya'ken !kei-ya'.

!gou i' #huai ha'hi

!ho'ari !huonni !i'ke'.



BUSHMAN TEXT



- Diacritics cannot be represented using Unicode
- No one left that speaks the |xam language!
- Over 137 different diacritics (more still being found)

a a' a'' a''' a''''





ENCODING



- Bushman text cannot be encoded using Unicode
- Latex IPA package contains diacritics
 - Allows for custom macros to be created
- Stacked, nested, multiple characters
- $\backslash\text{uline}\{a\} \rightarrow \overset{\cdot}{\underset{\cdot}{a}}$
- $\backslash\text{xbelow}\{\backslash\text{uline}\{a\}\} \rightarrow \overset{\cdot}{\underset{x}{\underset{\cdot}{a}}}$
- $\backslash\text{xbelow}\{aa\} \rightarrow \underset{x}{aa}$





ENCODING



!gaù ě #kuàni kàu¹kí

!ga\dialine{u} \twodotsu{e}
 \texthash{}ku\barblinet{a}n
 \ybelow{k}a\uline{u}k\u{i}

!gaù ě #kuà_yn kàu¹kí

!kù ě yàkèñ !kui-yà¹.

!k\uline{u} \twodotsu{i}
 y\dialine{a}ke\u{n}
 !ku\uline{i}-y\uline{a} .

!kù ě yàkèñ !kui¹-yà¹ .

!kò¹āñ k¹uonni¹ !ùhè¹.

!k\barbelow{\dialine{o}}
 \circbtwodotsa{a}\onedot{n}
 \xcbelow{k}\uline{uo}nn\u{i}
 !\uu{u}h\uline{e} .

!kò_oāñ k_c¹uonni¹ !ù¹hè¹ .





XÒÄ'XÒÄ - “TO WRITE”



- An AJAX tool to create a Bushman corpus
- Automatic algorithms
- User input
- Preprocessing
- Line and word segmentation
- Transcription
- Job and user management





LINE SEGMENTATION

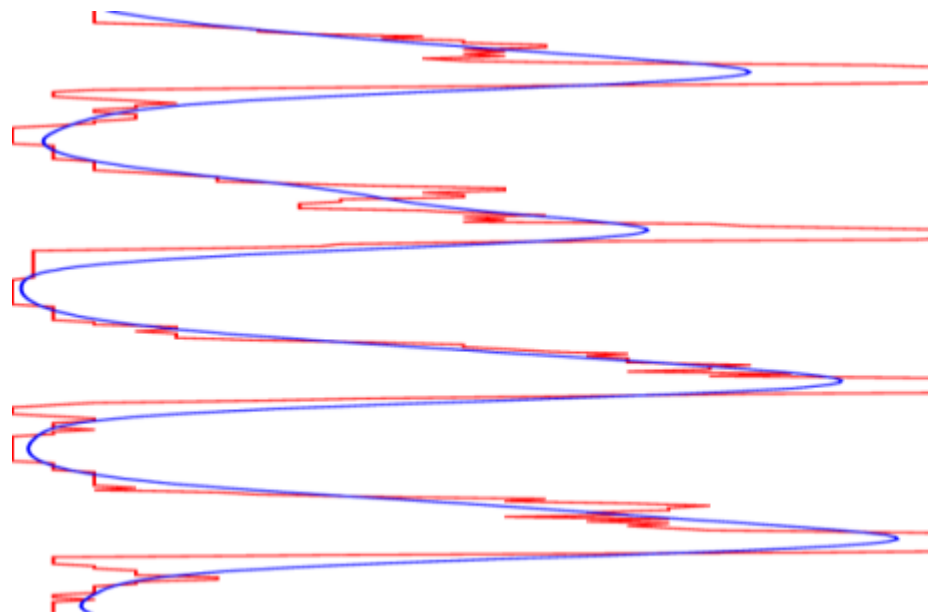


- Projection profile-based line segmentation
- Count foreground-background transitions for each row
- Minima suggest space between lines
 - Could represent space between base character and diacritics
 - Gaussian smoothing of projection profile





LINE SEGMENTATION



$t\bar{a}^1, \#kha\ddot{u}khe\ddot{n} + k\ddot{u}$
 $!na\ddot{u}n - k\ddot{h}\ddot{o} + k\ddot{h}\ddot{i}i.$
 $!k\ddot{u} - k\ddot{h}\ddot{o} - k\ddot{h}\ddot{e}n + k\ddot{h}\ddot{a}k\ddot{h}\ddot{o}$
 $ka, \#kha\ddot{u}khe\ddot{n} k\ddot{h}\ddot{a}n$
 $+ k\ddot{u} !k\ddot{a}nna, k\ddot{h}\ddot{i}i + k\ddot{h}\ddot{i}i.$



LINE SEGMENTATION



Menu

Home
New Segmentation Job
New Transcription Job
Logout



4785
he' ase' + kha'
thri, he' ase' a'
he' ta' ! kha' kha'
a', kha'. Ty' herre-
ten + ne' kha' - kha'

Info

The algorithm has identified the following line segmentation points. Please check to see if they are right.

You can **move** any segmentation point by clicking on it and dragging it to a new location.

You can **add** a new segmentation by clicking the add button.

You can **delete** a segmentation point by right-clicking on it.

When you are done and if you are satisfied with all of the segmentation points, click "Segment Lines!" to continue.

Change the candidates

If you are getting bad results you can try improve them by selecting an intensity level below and clicking 'Try again!'

- ☐ Lowest
☐ Low
☐ Medium
☐ High
☐ Very High

Try Again!



Digital Libraries Laboratory, University of Cape Town



WORD SEGMENTATION



- Line slant is automatically corrected
- Connected components in text lines are identified
- Distances between adjacent components are calculated
- Distances above threshold separate words





WORD SEGMENTATION



Menu

Home
New Segmentation Job
New Transcription Job
Logout



Hoi, he se a

New Segmentation Point

he ta ! Maiken

New Segmentation Point

a, ha. Tyherre

New Segmentation Point

Segment Words

Info

The algorithm has identified the following word segmentation points. Please check to see if they are right.

You can **move** any segmentation point by clicking on it and dragging it to a new location.

You can **add** a new segmentation by clicking the add button.

You can **delete** a segmentation point by right-clicking on it.

When you are done and if you are satisfied with all of the segmentation points, click "Segment Words!" to continue.



Digital Libraries Laboratory, University of Cape Town





CORPUS CREATION WORKSHOPS



- Workshop held to create Bushman corpus
- 29 data capturers recruited
- 900 pages from 2 authors randomly selected
- 729 pages were segmented into lines and words
- 1547 text lines were transcribed
- 452 text lines could not be transcribed
 - Interface didn't support characters, noise, English





CORPUS CREATION WORKSHOPS



-
- Quality and efficiency of data capturers evaluated
 - 5 data capturers asked to return
 - 1700 more line recruited
 - More efficient and potentially fewer errors





CORPUS CREATION WORKSHOP



Digital Libraries Laboratory, University of Cape Town



USER CONTRIBUTIONS



Segmentation Jobs		Transcription Jobs	
Jobs	Number of Users	Jobs	Number of Users
0-9	0	0-9	0
10-19	5	10-19	1
20-29	9	20-29	3
30-39	7	30-39	4
40-49	5	40-49	9
50-59	0	50-59	4
60-69	0	60-69	6
70-80	0	70-80	2





DATA QUALITY

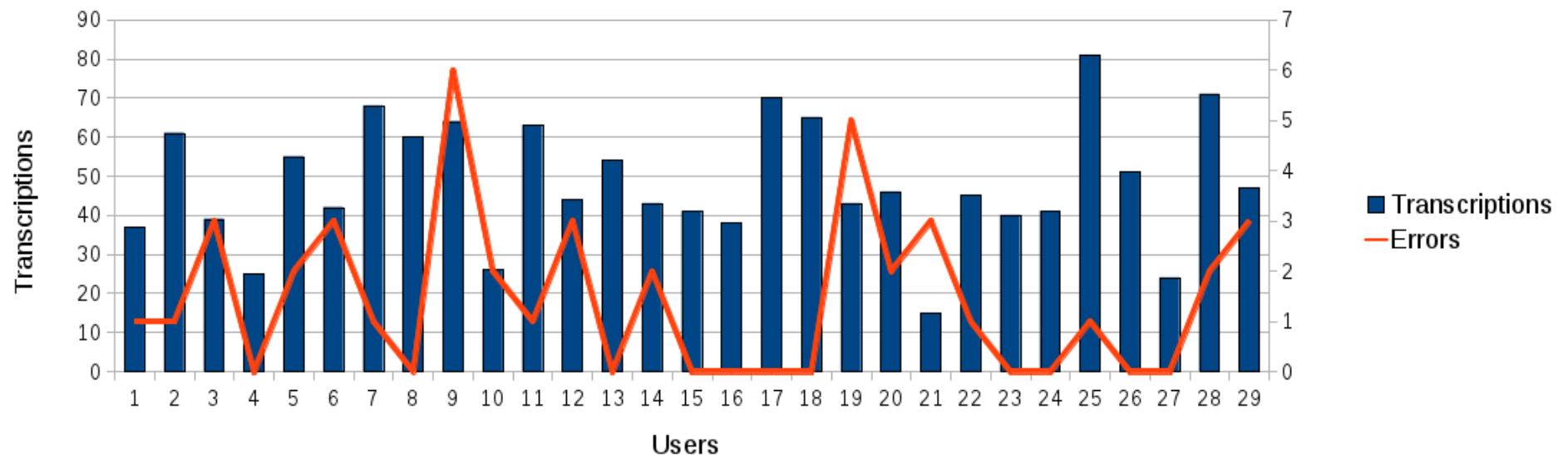


- Quality represented by accuracy and correctness of transcriptions
- Useful in planning for follow on workshops
- Random transcriptions by each user evaluated by research assistant
 - Wrong diacritics, characters, etc.
- Average of 0.48 errors per text line
 - Acceptable for lay persons?





EFFICIENCY VS QUALITY



Digital Libraries Laboratory, University of Cape Town



CONCLUSIONS



-
- Creation of corpora for historical texts is often difficult due to complexities of script
 - Semi-automatic tool allowed for more efficient and less expensive creation of corpus
 - Currently being used in handwriting recognition study
 - Applicable to other historical collections





THANK YOU



Questions?



Digital Libraries Laboratory, University of Cape Town