Printed-Source Analysis for PanLex

Jonathan Pool The Long Now Foundation





First PanLex Summer Internship 29 July 02013

Summary

- Printed sources
- PanLex collection at Internet Archive
- Analytical tasks
- Tools
- Opportunities



Printed formats



Printed sources in the source archive

plx=# sele count		
3993	<pre>plx=# select count (ap) f (select ap from ap except</pre>	rom select
274 printed- only	af.ap from af, fm where f af.fm and (tt like 'pdf-% like 'prt%')) as tbl; count 3719	m.fm = 'or tt

Printed sources in the world

			OPEN			
SUBJECTS						
AUTHORS	ADD A BOOK		LIBRARY			
LISTS	RECENTLY	HELP	One web page for every book.			
Search Results 66 773 hits = Relevance Most Editions First Published Most Recent						
subject:d	lictionaries		S			

Printed sources in the world

Polyglot Dictionaries 1 140 works

Dictionaries, Polyglot 158 works



Search Results

Results: 1 through 36 of 36 (0.001 secs)

 <u>The Chimwiini lexicon exemplified</u> - Kisseberth, Charles W Includes bibliographical references (p. xxxiv) Keywords: <u>Mwini dialect</u>; <u>Swahili language</u> Downloads: 3
 <u>A classified vocabulary of the Punu language</u> - 880-01 Yukawa, Yasutoshi, 1941- Includes index <u>Keywords: Punu language</u>; <u>Punu language</u>; <u>English language</u>; <u>French language</u>; <u>Punu (African people)</u>; <u>Punu</u> <u>Downloads</u>: 1
 <u>Degema-English dictionary with English index</u> - Kari, Ethelbert E Includes bibliographical references (p. xliii-xliv) and index Keywords: <u>Degema language</u>; <u>English language</u>; <u>Lengua degema</u>; <u>Lengua inglesa</u>; <u>Lengua degema</u>; <u>Lengua inglesa</u> Downloads: 2
 <u>EG-Woerterbuch mathematischer Begriffe = EK-vortaro de mathematikaj terminoj</u> - Higers-Yashovardhan, R. edt Keywords: <u>Mathematik</u>; <u>Europäische Gemeinschaft</u>; <u>Mathematik</u>; <u>Wissenschaftliche Bibliothek</u> Downloads: 1
 Gendal Ulgurugo shojiten = Concise modern Uyghur-Japanese dictionary - 880-01 Sugawara, Jun, 1966- Parallel title also in Ulghur Keywords: Ulghur language; Ulgurugo Downloads: 4
○ 『ヒンディー語動詞基礎語彙集』 = Hindi - Japanese dictionary of selected verbs with illustrative sentences - 町田和彦/Kazuhiko Machida Downloads: 4
 A Jita vocabulary - Kagaya, Ryohei Keywords: Jita language; Jita language; English language; Japanese language; Jita; Engels; Japans Downloads: 2
 <u>Guyish-i Gilaki-i Lahijan : vazhahnamah va parahi vizhagiha-yi avayi va sakhtavazhahi</u> - Jahangiri, Nader jild-i 1. Vizhagiha-yi avayi va sakhtavazhahi jild-i 2. Vazhagan-i Alif ta Sin jild-i 3. Vazhagan-i Shin ta Ya Keywords: <u>Gilaki language</u>; <u>Gilaki language</u>; <u>Gilaki;</u> <u>Gilaki language</u> Downloads: 1
 Gendal Chibettogo doshi jiten : Rasa hogen = A verb dictionary of the modern spoken Tibetan of Lhasa, Tibetan-Japanese - 880-01 Hoshi, Izumi Dictionary entries in Tibetan script with romanization Keywords: <u>Tibetan language</u> Downloads: 1
 Kainanto hogen kiso golshu - 880-01 Nakajima, Motoki, 1942- Includes bibliographical references (p. 8) and index Keywords: <u>Chinese language</u> Downloads: 1
 ○ AA諸言語教育基本語彙表(入門期の学習に必要な基礎語彙600項目試案) - 梅田博之 Downloads: 4
 <u>The Pyen (or Phen) language : its classified lexicon</u> - 880-01 Shintani, Tadahiko, 1946- Downloads: 3

Gendai Uigurugo shojiten = Concise modern Uyghur-Japanese dictionary (2009)

Title (alternate script): = Concise modern Uyghur-Japanese dictionary Author: 880-01 Sugawara, Jun, 1966-Subject: Uighur language; Uigurugo Publisher: Fuchu-shi : Tokyo Gaikokugo Daigaku Ajia Afurika Gengo Bunka Kenkyujo Language: Japanese; Uighur Digitizing sponsor: Jonathan Pool Book contributor: The Long Now Foundation, PanLex Project Collection: lendinglibrary; browserlending; panlex; longnow; americana Notes: some content may be lost due to the binding of the book.

Full catalog record: MARCXML

mathis book has an editable web page on Open Library.

Description

Parallel title also in Uighur

Includes bibliographical references (p. xxii-xxiv)

Selected metadata

Isbn:	9784863370265					
Lccn:	2010428072					
Page-progression:	Ir					
Scanningcenter:	sanfrancisco					
Mediatype:	texts					
Shiptracking:	y0599					
Identifier:	gendaiuigurugosh008800					
Scanner:	scribe5.sanfrancisco.archive.org					
Ppi:	500					
Camera:	Canon EOS 5D Mark II					
Operator:	associate-lan-zhu@					
Scandate:	20130408202946					
Republisher:	associate-karina-martinez@					
Imagecount:	788					
Identifier-access:	http://archive.org/details/gendaiuigurugosh008800					
Identifier-ark:	ark:/13960/t8ff5b645					
Bookplateleaf:	0004					
Ocr:	ABBYY FineReader 8.0					
Sponsordate:	20130430					

PanLex 3.4

source — 3317 — jpn-uig:菅原

number	3317
PanLex — beginning	2011-04-25
name	jpn-uig:菅原
World Wide Web	http://www.aa.tufs.ac.jp/en/publications/lexicon
ISBN	9784863370265
author	菅原純(編)
title	『現代ウイグル語小辞典』 = Concise modern Uyghur-Japanese dictionary
publisher	アジア・アフリカ言語文化研究所
year	2009
person — good	6
person — number	
fact other	アジア・アフリカ基礎語彙集/Asian and African Lexicon, 53 (B026)
permission - kind	na
right text	
right - person - name	
permission - email - address	

good - number - edit - necessar	y — 1/0 0
file —	difficult 8
file — submit — necessar	y — 1/0 1
expression — edit — necessar	y — 1/0 0
expression - edit - done - language - (aaa,b	bbb,ccc)
file —	address jpn-uig-菅原
fact	— other book donated by ILCAA

translation - count

language

jpn-000	日本語	
uig-000	ئۇيغۇرچە	
uig-001	Uyghurche	

file — kind

prt@Plx

edit — person

file — difficult	8
file — submit — necessary — 1/0	1
expression — edit — necessary — $1/0$	0
expression — edit — done — language	
— (aaa,bbb,ccc)	
file — address	jpn-uig-菅原

PanLex 3.4

[PanLem (usred6w)] [9.1.9]

you - testintern

beginning

back

source — edit — permission

	change	source	language —	file — kind	difficult	title
			necessary			
Г						

	1				
±	nan-eng:Nakajima	nan-000	pdf-img	8	『海南島方言基礎語彙集』
±	eng-jpn-rus-niv:丹 菊逸治	niv-000; jpn-002	pdf-img	8	ニヴフ語サハリン方言基礎語彙 集(ノグリキ周辺地域) = Basic vocabulary of the Sakhalin dialect of Nivkh language : Nogliki dialect
±	nsz-eng:Kroeber	nsz-000	pdf-img	4	The Valley Nisenan
±	ood-eng:Allison	ood-000	pdf-img	4	O'otham Ñiokĭ Haichu A:ga
±	eng-mya-shn- pce:Shin	pce-002; pce-001; mya-002; shn-000; pce-003; pce-004; pce-005	pdf-img	7	The Palaung Language: Comparative Lexicon of its Southern Dialects (I)

Same as for digital text, plus:

- Page imaging
- Spatial structural analysis
- Text identification and sequencing
- Script identification
- Font identification
- Face identification
- Character segmentation and sequencing
- Character identification

Page imaging

- Produces image file(s)
- Internet Archive produces JPG2000, PDF/A, Djvu, and Ebook files (in color)
- PanLex collection: format changes from "prt@Plx" to "PDF-img"



Page imaging

• Internet Archive resolution



Spatial structural analysis

1	HE' Yan	72 10				imantru
1	Y Kyan	0. 001.				ウシ科の家
	npf.	pf.		inıp.		寝具を干す
	নন্তু ন'	지劧도지.		র্ন্টুদশ		RT.~ Thor
	`kyan	`kyan		`kyan		
	pres.	fut.	pt.	inın.		与具を乾か
	r ¥=.			¥	,	•9`~ ~sha`l
	ι ψ _γ	A T	NULA	U LA	1	ノコを干す
	1. (~を)伸は	টের নানামা	\sim -kang	pa `kyan ½	こを伸 しんしょう しんしょ しんしょ	
	ばす. ホニ・マーリ	\sim `kanglaa	`kyan ₹	足を伸ば	इ. हे	- Model -
	\sim -ce `kyan f	舌を出す. つ	<u>द्र</u> ेंद्र:यॅ:~	′drongko ′	`kyan	지ᇌ [·] ·ku
	真っ直ぐに伸ば	র্টক. শাত্রমার	মশা~`yä	iälaa `kyar	n右手	npf.
	をあげる. 94	ママー~ ´lakpa	`kyan ①	手を伸ばす	す. ②	고왔
	けちる.③手	を出す. 殴る	5.			-ku

Text identification and sequencing



Script identification

CEPA

 (\dot{H})

S cepa СЕРА (лат. *Sulfur*), S, 32,066. Химический с четырех стабильных I ³⁴S (4,16%) и ³⁶S (0, 0,170 нм (координаци 6) и иона S⁶⁺ 0,026 н ионизации нейтральн

23,35, 34,8, 47,3, 72,5 и 88,0 эВ. Сера Менделеева, в 3-м периоде, и принад электронного слоя 3s²3p⁴. Наиболее (валентности соответственно II, IV и Сера относится к числу неметаллов.

Center for European Policy Analysis

Font identification

0123456789 Illinois 0123456789 Illinois 0123456789 Illinois

Face identification

0123456789 Illinois

0123456789 Illinois

Character segmentation and sequencing



Character identification



Image-to-text conversion ("optical character recognition") tools: competitors

- Nicomsoft
- Tesseract
- OCRopus
- Abby FineReader
- Adobe Acrobat
- OmniPage
- Readiris
- GOCR
- LEADTOOLS
- Many others

Image-to-text conversion ("optical character recognition") tools: properties

- Commercial versus open-source
- Script-specific versus multiscriptal
- Language-specific versus multilingual
- Diaglottal versus synglottal (language-mixing)
- Immutable versus trainable
- Isolated versus integratable
- Fossilized versus actively developed

Prior research

- Christa Mabee, "<u>Optical Character Recognition in Multilingual Text: A Brief</u> <u>Survey</u>", 02012
- Christa Mabee, "<u>Report on Internship Project: Optical Character Recognition in</u> <u>Multilingual Text</u>", 02012
- Marcin Heliński, Miłosz Kmieciak, Tomasz Parkoła, "<u>Report on the comparison of</u> <u>Tesseract and ABBYY FineReader OCR engines</u>", 02012 (<u>PanLex copy</u>)
- Ray Smith, Daria Antonova, and Dar-Shyang Lee, <u>"Adapting the Tesseract Open</u> <u>Source OCR Engine for Multilingual OCR</u>", 02009
- Burcu Karagol-Ayan, "<u>Resource Generation from Structured Documents for Low-</u> <u>density Languages</u>", 02007
- Tapas Kanungo and Song Mao, "<u>Stochastic Language Models for Style-Directed</u> <u>Layout Analysis of Document Images</u>", 02003

Prior research: conclusions

- Page imaging is hard
- Tool training is necessary but expensive
- Non-dictionary corpus integration has not been tried
- PanLex requires multiscriptal, multilingual, synglottal, trainable, integratable, actively developed tools
- Few (≈ 2) tools even aspire to all of these properties
- Most likely candidates now: ABBYY FineReader 11 Professional Edition, ABBYY FineReader Engine, Tesseract

ABBYY FineReader 11 Professional Edition

- Application
- Commercial
- SKU FRPFWA11E
- Windows 8 / 7 / Vista / XP; Windows Server 2012 / 2008 / 2008 R2 / 2003
- Express edition (Windows and Macintosh) and web service are more limited
- Price per license: list \$300, street \$270+, direct (through 31 July) \$180 (or \$120 for "upgrade")

ABBYY FineReader Engine

- SDK
- Commercial
- Version 10: Windows 7 / Vista / XP / 2000; Windows Server 2003
- Version 9: Linux (GCC 2.95-4)
- Version 8: Macintosh (OS X 10.4–10.8)
- Prices not published

Tesseract

- SDK and command-line application
- Open-source (Apache License)
- Partly sponsored by Google
- Version 3.02: Linux, OS X
- Additional OSs: Windows with VC++ or CygWin, iOS, Android
- Used by about 40 other tools and projects

Script identification

enersen heren h	AFRP 11	AFRE 10	AFRE 9	AFRE 8	Tesseract
Latin	۲		۲	۲	
Cyrillic	۲	۲	۲	۲	
Greek	۲		۲	۲	
Armenian	۲		۲	۲	
Hebrew	۲			۲	
Han		۲	۲		
Hiragana/Katakana	۲	۲	۲		
Hangul	۲	۲	۲		
Thai	۲	۲	۲		
Arabic	۲	()			
Devanagari					
Telugu					
Tamil					
Malayalam					
Kannada					
Bengali					
Cherokee					

Quality comparison

- Heliński et al.:
 - FineReader was better on average than Tesseract, but worse on some tasks
 - Test was entirely on Polish (thus Latin-script) text
- Mabee:
 - FineReader was better on average than Tesseract
 - Test was on Latin, Thai, Arabic, and Tibetan scripts

Training comparison

- Both FineReader and Tesseract can (in principle) be trained, but:
 - Bugs prevented Mabee from training FineReader
 - Mabee found Tesseract training tedious and useless
- Approximate training costs for 1 font per Heliński et al.:
 - FineReader by ABBYY: \$40,000
 - FineReader in-house: 25 person-days
 - Tesseract in-house: 8 person-days

Extensibility comparison

- Tesseract operates on more scripts than FineReader
- But both omit Georgian, Sinhala, Lao, Tibetan, Burmese, Khmer, Ethiopic, etc.
- So both need to be extended for PanLex
- Both permit defining new languages and providing new dictionaries

Integratability comparison

- Desiderata for PanLex:
 - Dynamic integration of expression corpus (ex table)
 - Dynamic integration of approved characters (cu and cp tables)
- ABBYY FineReader 11 Professional Edition: no
- ABBYY FineReader Engine: maybe
- Tesseract: probably

Performance on an easy source

• FineReader

S

ruvengo hatred ruzha noise -rwa v. fight -rwara be ill; sick rwendo journey rwiyo song

sachigaro chairperson

ruvengo hatred
ruzha noise
-rwa v. fight
-rwara be ill; sick
rwendo journey
rwiyo song

sachigaro chairperson

Performance on an easy source

• Tesseract

S

ruvengo hatred ruzha noise -rwa v. fight -rwara be ill; sick rwendo journey rwiyo song

sachigaro chairperson

ruvengo hatred ruzha noise -rwa v. fight -rwara be ill; sick rwendo journey rwiyo song S sachigaro chairperson

Performance on a more complex source

• FineReader

 kạri² [karí] vi decree good things
 kel [kél] vi

 kạsany [kasáŋ] vi cough USAGE May
 yọ 'There

 kari 2 [kari] vt decree good things

 kasany [kasan] vi cough USAGE May

Performance on a more complex source

• Tesseract

 kạri² [karí] vt decree good things
 kel [kél] vt

 kạsany [kasáŋ] vi cough USAGE May
 yọ 'There

 kariz [kari] vt decree good things

 kasany [kasap] vi cough USAGE May

• FineReader	npf.	pf.		imp.		
	নস্ত্রন.	고월다.		କ୍ଷିଦା, କ୍ରିଦା.		
	¯kyää	-kyää		-kyää, -kyöö		
	pres.	fut.	pt.	imp.		
	[ફ્રેવ	고월지.	ସକ୍ରିଦା.	સુંત.]	



 Fine ("Fin HTN outp 	Reader nish", ML ut)	npf ק	्या vää pres.		pf. 고괽지 ⁻ kyää fut. 지괽지	pt. 지원대	imp. 뒷작, 풋작 Tkyää, 다 imp. 풋작	cyöö]
<i>npf</i> . qg^' "kyää	P 9 "1	9 <i>f-</i> spr kyää		imp go ''ky). ,-, q- /ää, "kyöö	ö		
[ar	qgof	qgq	-	aj-]			

• Tesseract	esseract npf. সম্ভূন		<i>pf.</i> 지원다		imp. સુત્ર", સુંત્ર"	
	⁻kyää		⁻kyää		⁻kyää, ⁻kyöö	
		pres.	fut.	pt.	imp.	
	[શ્રુત.	ସକ୍ରିଦା.	ସକ୍ତ୍ରିଦା.	સુંત.]

```
npf. pf. Imp.
magma" r1504" gm", §m'
'kyééi 'kyéi5 'kyéiéi, 'ky66
pres. fut. pt. mzp.
[ gar agar qgw"
\F
```

Performance on a very complex source

 Tesseract ("Finnish") 	np - k	f. yar yää	pf. সম্ভূন্ম Tkyää		imp. 광대, 췴대 Tkvää, Tkvöö	
	[pres. শ্রুনা	fut. 지원지	pt. 지원지	imp. সুন]

```
npf. pf. ımp.
mgflı' rıgfıı' gfır, šfıı'
-kyää -kyää -kyää, -kyöö
pres. fııl. pt. mıp.
[ àfıı' ngfıı' fiıgw'
`f
```

K

Performance on a very complex source

• Tesseract	np	f.	pf.		imp.	
("Finnish",	নস্ত্রন'		নম্ভূন	•••	କ୍ରିୟ, କ୍ରିୟ,	
HIML	k	yaa	-kya	a	-kyaa, -k	cyoo
output)		pres.	fut.	pt.	imp.	
	[ਉਪ.	고원다.	ଅକ୍ତିମ.	রীন.]

npf. pf. 1mp. mgfl1' rıgfı1' gfır, šfı1'

-kyää -kyää -kyää, -kyöö

pres. fiil. pt. mip.

[àfii' ngfii' fiigw'

`f

Opportunities

PanLex and image-to-text conversion tools

- Benefits to PanLex from tool extensions
- Benefits from PanLex for tool extensions
- Which tool(s) to use and support?
- Multiple tools?
 - Advantages: federation, complementarity, independence
 - Disadvantages: expense, forfeiture of expertise

Opportunities

Partnerships for tool development

- Tool developers
- Possessors of lexical data
- Tool consumers
 - Internet Archive
 - HathiTrust
 - Google
 - Project Gutenberg