

Printed-Source Analysis for PanLex

Jonathan Pool
The Long Now Foundation

PanLex



<http://panlex.org>

First PanLex Summer Internship

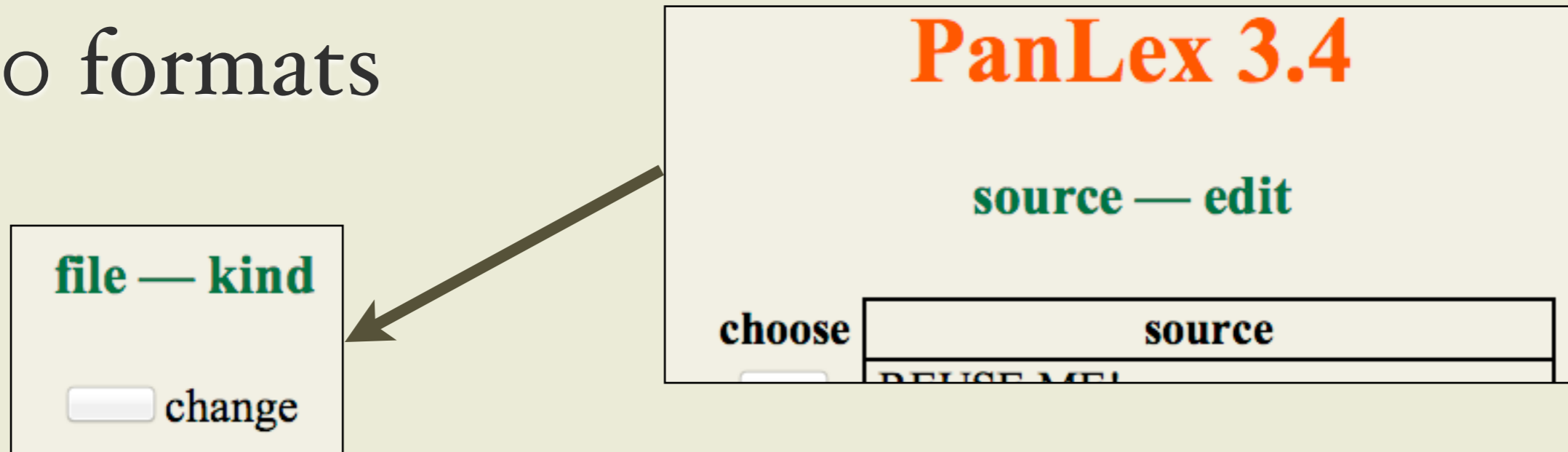
29 July 02013

Summary

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

Printed sources

70 formats





PanLex 3.4

source — aak-eng:Speece — file — kind

GDN	<input type="checkbox"/>	PCS
Geiriadur	<input type="checkbox"/>	NIL; #<DICTIONARY:LOOKUP-PROXY; []
JMdict	<input type="checkbox"/>	<!DOCTYPE JMdict; <entry>; <sense>; <gloss xml:lang=
LL-LIFT	<input type="checkbox"/>	linguist-lift-schemas/linguist-lift
LTRC	<input type="checkbox"/>	"_७०ड", "अव्यय", "avy", "_के_जैसा"

Printed sources

Printed formats

pdf	<input type="checkbox"/>	
pdf-img	<input type="checkbox"/>	
pdf-lock/encrypt	<input type="checkbox"/>	
png	<input type="checkbox"/>	
po	<input type="checkbox"/>	
popdict	<input type="checkbox"/>	
prc	<input type="checkbox"/>	
prt	<input type="checkbox"/>	 @ ⊕
prt@Plx	<input type="checkbox"/>	
ps	<input type="checkbox"/>	

Printed sources

Printed sources in the source archive

```
plx=# select count (ap) from ap;  
count
```

```
-----  
3993
```

274
printed-
only

```
plx=# select count (ap) from  
(select ap from ap except select  
af.ap from af, fm where fm.fm =  
af.fm and (tt like 'pdf-%' or tt  
like 'prt%')) as tbl;
```

```
count  
-----  
3719
```

Printed sources

Printed sources in the world

The screenshot shows the Open Library website interface. At the top right, there is a logo with a blue sign that says "OPEN" hanging from a string, with the word "LIBRARY" in a large serif font below it. Underneath the logo is the tagline "One web page for every book." On the left side, there is a navigation menu with links for "SUBJECTS", "AUTHORS", "LISTS", "ADD A BOOK", "RECENTLY", and "HELP". The main content area displays "Search Results" in a large teal font. Below this, it shows "66 773 hits" in green, followed by a hamburger menu icon and several sorting options: "Relevance", "Most Editions", "First Published", and "Most Recent". At the bottom, there is a search bar containing the text "subject:dictionaries" and a "Search" button partially visible on the right.

SUBJECTS

AUTHORS

LISTS

ADD A BOOK

RECENTLY

HELP

OPEN

LIBRARY

One web page for every book.

Search Results

66 773 hits ☰ [Relevance](#) | [Most Editions](#) | [First Published](#) | [Most Recent](#)

subject:dictionaries

Search

Printed sources

Printed sources in the world

Polyglot Dictionaries 1 140 works

Dictionaries, Polyglot 158 works

Internet Archive PanLex collection



The screenshot shows the Internet Archive website interface. At the top left is the logo with the text "INTERNET ARCHIVE" and a classical building icon. To the right are navigation links: "Web", "Video", "Texts", and "Audio". Below this is a secondary navigation bar with "Home", "donate", "Forums", "FAQs", and "Collections". A search bar contains the text "Search: collection:panlex". Below the search bar, the heading "Search Results" is displayed in large bold font, followed by the text "Results: 1 through 36 of 36 (0.001 secs)".

INTERNET ARCHIVE

Web Video Texts Audio

Home donate | Forums | FAQs | Collections

Search: collection:"panlex"

Search Results

Results: 1 through 36 of 36 (0.001 secs)

Internet Archive PanLex collection

- [A comparative vocabulary of Formosan languages and dialects](#) - Ogawa, Naoyoshi
Keywords: [Austronesische talen](#)
Downloads: 2
- [The Chimwiini lexicon exemplified](#) - Kisseberth, Charles W
Includes bibliographical references (p. xxxiv)
Keywords: [Mwini dialect](#); [Swahili language](#)
Downloads: 3
- [A classified vocabulary of the Punu language](#) - 880-01 Yukawa, Yasutoshi, 1941-
Includes index
Keywords: [Punu language](#); [Punu language](#); [Punu language](#); [English language](#); [French language](#); [Punu \(African people\)](#); [Punu](#)
Downloads: 1
- [Degema-English dictionary with English index](#) - Kari, Ethelbert E
Includes bibliographical references (p. xliii-xliv) and index
Keywords: [Degema language](#); [English language](#); [Lengua degema](#); [Lengua inglesa](#); [Lengua degema](#); [Lengua inglesa](#)
Downloads: 2
- [EG-Woerterbuch mathematischer Begriffe = EK-vortaro de matematikaj terminoj](#) - Hilgers-Yashovardhan, R. edt
Keywords: [Mathematik](#); [Europäische Gemeinschaft](#); [Mathematik](#); [Wissenschaftliche Bibliothek](#)
Downloads: 1
- [Gendai Uigurugo shojiten = Concise modern Uyghur-Japanese dictionary](#) - 880-01 Sugawara, Jun, 1966-
Parallel title also in [Uighur](#)
Keywords: [Uighur language](#); [Uiguruqo](#)
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](#); [Jita language](#); [English language](#); [Japanese language](#); [Jita](#); [Engels](#); [Japans](#)
Downloads: 2
- [Guyish-i Gilaki-i Lahijan : vazhahnamah va parahi vizhagiha-yi avayi va sakhtavazhahi](#) - Jahangiri, Nader
jild-i 1. Vzhagiha-yi avayi va sakhtavazhahi -- jild-i 2. Vzhagan-i Alif ta Sin -- jild-i 3. Vzhagan-i Shin ta Ya
Keywords: [Gilaki language](#); [Gilaki language](#); [Gilaki](#); [Gilaki language](#)
Downloads: 1
- [Gendai 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: [Tibetan language](#)
Downloads: 1
- [Kainanto hogen kiso goishu](#) - 880-01 Nakajima, Motoki, 1942-
Includes bibliographical references (p. 8) and index
Keywords: [Chinese language](#)
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
- [The Pyen \(or Phen\) language : its classified lexicon](#) - 880-01 Shintani, Tadahiko, 1946-
Downloads: 3

Internet Archive PanLex collection

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](#)

 This book has an [editable web page](#) on [Open Library](#).

Description

Parallel title also in Uighur

Includes bibliographical references (p. xxii-xxiv)

Internet Archive PanLex collection

Selected metadata

Isbn:	9784863370265
Lccn:	2010428072
Page-progression:	lr
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

Internet Archive PanLex collection

PanLex 3.4

source — 3317 — jpn-uir:菅原

number	3317
PanLex — beginning	2011-04-25
name	jpn-uir: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 — necessary — 1/0	0
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-uir-菅原
fact — other	book donated by ILCAA

translation — count

language

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

file — kind

edit — person

Internet Archive PanLex collection

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

Internet Archive PanLex collection

 stop

 beginning

 back

PanLex 3.4

[PanLem (usred6w)]
[9.1.9]

you — testintern

source — edit — permission

change	source	language — necessary	file — kind	difficult	title
<input type="checkbox"/> ±	nan-eng:Nakajima	nan-000	pdf-img	8	『海南島方言基礎語彙集』
<input type="checkbox"/> ±	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
<input type="checkbox"/> ±	nsz-eng:Kroeber	nsz-000	pdf-img	4	The Valley Nisenan
<input type="checkbox"/> ±	ood-eng:Allison	ood-000	pdf-img	4	O'otham Ńiokĩ Haichu A:ga
<input type="checkbox"/> ±	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)

Analytical tasks

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

Analytical tasks

Page imaging

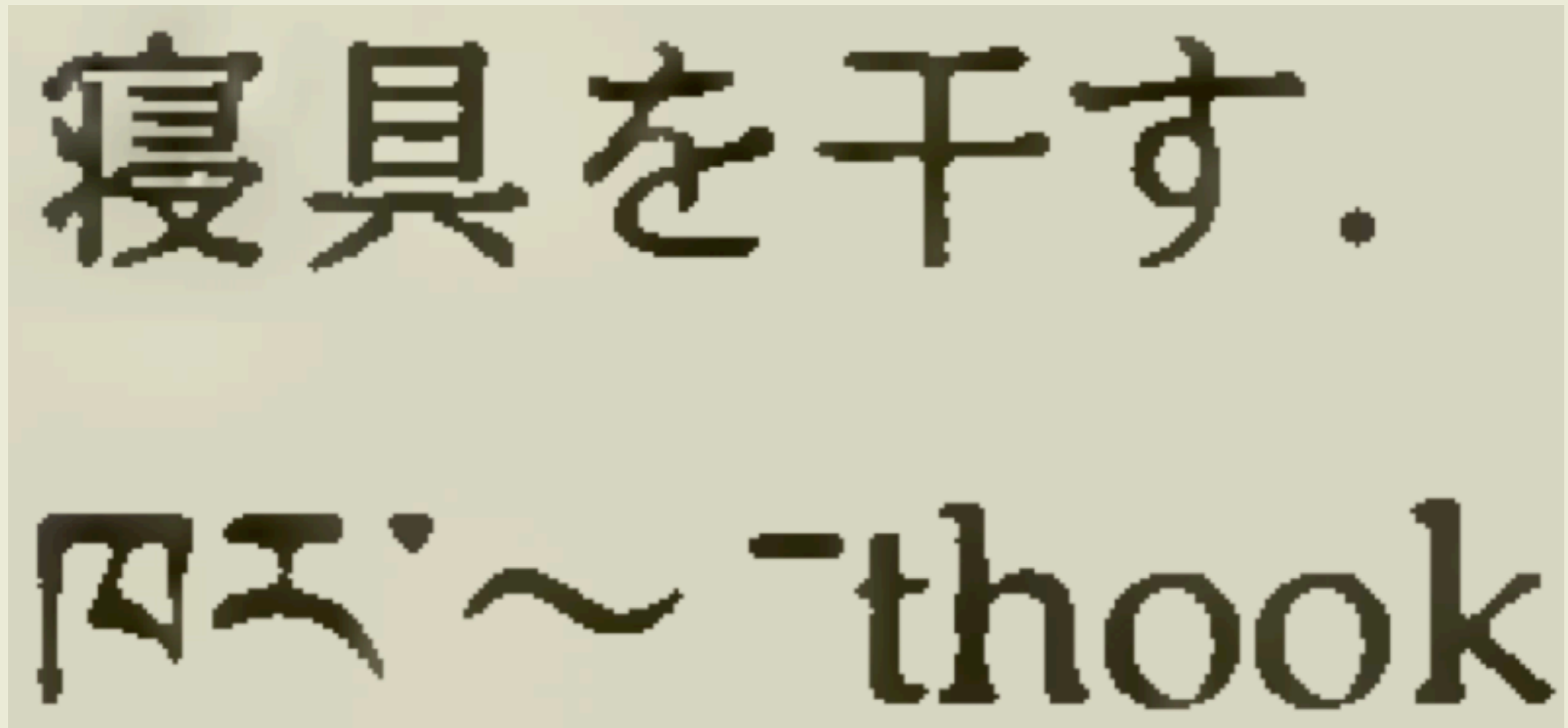
- Produces image file(s)
- Internet Archive produces JPG2000, PDF/A, Djvu, and Ebook files (in color)
- PanLex collection: format changes from “prt@Plx” to “PDF-img”



Analytical tasks

Page imaging

- Internet Archive resolution



Analytical tasks

Spatial structural analysis

བཀྱང་	`kyan	v. vol.							
<i>npf.</i>		<i>pf.</i>		<i>imp.</i>					
བཀྱང་		བཀྱངས་		ཀྱངས་					
`kyan		`kyan		`kyan					
<i>pres.</i>		<i>fut.</i>		<i>pt.</i>		<i>imp.</i>			
[ཀྱང་		བཀྱང་		བཀྱངས་		ཀྱངས་			

1. (~を) 伸ばす. ཀང་པ་~ `kangpa `kyan 足を伸ばす. ཀང་ལག་~ `kanglaa `kyan 手足を伸ばす. ལེ་~ `ce `kyan 舌を出す. འདྲོང་པོ་~ `drongko `kyan 真っ直ぐに伸ばす. གཡས་ལག་~ `yäälaa `kyan 右手をあげる. ལག་པ་~ `lakpa `kyan ① 手を伸ばす. ② けちる. ③ 手を出す. 殴る.

Umantra
ウシ科の家
寝具を干す.
ཁང་~ `thoo
写真を乾か
ག་~ `sha `k
ノコを干す.

བསྐྱེས་པ་
བསྐྱེ་ `ku
npf.
བསྐྱེ་
`ku

Analytical tasks

Text identification and sequencing

F

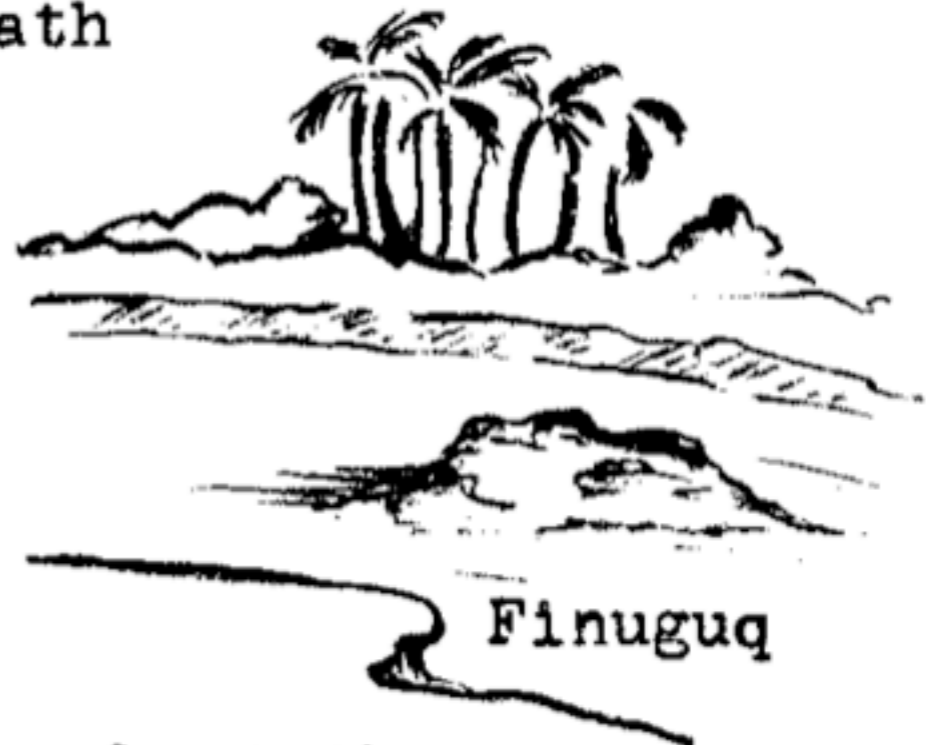
FETFET Whistle by drawing breath
over lower teeth.

FINUGUQ (pulô) Island (in
river).

FOFOY (tutubí) Dragonfly.

FOOG (sipol, sutsot, paswit)
Whistle.

FOOLOK (diwatà) Type of fairy about the size of a



Analytical tasks

Script identification

СЕРА

СЕРА



СЕРА (лат. *Sulfur*), S, 32,066. Химический элемент, один из четырех стабильных изотопов: ^{34}S (4,16%) и ^{36}S (0,04%), радиус атома 0,170 нм (координационное число 6) и иона S^{6+} 0,026 нм (координационное число 6). Энергии ионизации нейтрального атома: 23,35, 34,8, 47,3, 72,5 и 88,0 эВ. Сера — неметалл Менделеева, в 3-м периоде, и принадлежит к группе 16. Валентный электронный слой $3s^2 3p^4$. Наиболее распространены валентности II, IV и VI. Сера относится к числу неметаллов.

СЕРА Center for European Policy Analysis

Analytical tasks

Font identification

0123456789 Illinois

0123456789 Illinois

0123456789 Illinois

Analytical tasks

Face identification

0123456789 Illinois

0123456789 Illinois

Analytical tasks

Character segmentation and sequencing

تسليم الحقائق

m m

हिन्दी

Analytical tasks

Character identification

1 1

one

ell

Tools

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

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

Tools

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

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

Tools

Prior research

- Christa Mabee, “Optical Character Recognition in Multilingual Text: A Brief Survey”, 02012
- Christa Mabee, “Report on Internship Project: Optical Character Recognition in Multilingual Text”, 02012
- Marcin Heliński, Miłosz Kmieciak, Tomasz Parkoła, “Report on the comparison of Tesseract and ABBYY FineReader OCR engines”, 02012 (PanLex copy)
- Ray Smith, Daria Antonova, and Dar-Shyang Lee, “Adapting the Tesseract Open Source OCR Engine for Multilingual OCR”, 02009

Tools

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 ($\cong 2$) tools even aspire to all of these properties
- Most likely candidates now: ABBYY FineReader II Professional Edition, ABBYY FineReader Engine, Tesseract

Tools

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”)

Tools

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

Tools

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

Tools

Script identification

	<i>AFRP 11</i>	<i>AFRE 10</i>	<i>AFRE 9</i>	<i>AFRE 8</i>	<i>Tesseract</i>
<i>Latin</i>	●	●	●	●	●
<i>Cyrillic</i>	●	●	●	●	●
<i>Greek</i>	●	●	●	●	●
<i>Armenian</i>	●	●	●	●	
<i>Hebrew</i>	●	●		●	●
<i>Han</i>	●	●	●		●
<i>Hiragana/Katakana</i>	●	●	●		●
<i>Hangul</i>	●	●	●		●
<i>Thai</i>	●	●	●		●
<i>Arabic</i>	●	(●)			●
<i>Devanagari</i>					●
<i>Telugu</i>					●
<i>Tamil</i>					●
<i>Malayalam</i>					●
<i>Kannada</i>					●
<i>Bengali</i>					●
<i>Cherokee</i>					●

Tools

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

Tools

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

Tools

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

Tools

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

Tools

Performance on an easy source

- FineReader

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

S

sachigaro chairperson



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

sachigaro chairperson

Tools

Performance on an easy source

- Tesseract

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

S

sachigaro chairperson



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

Tools

Performance on a more complex source

- FineReader

kari² [karí] vt decree good things kel [kél] vt
kasany [kasán] vi cough USAGE May yɔ 'There



kari 2 [kari] vt decree good things
kasany [kasan] vi cough USAGE May

Tools

Performance on a more complex source

- Tesseract

kari² [karí] vt decree good things kel [kél] vt
kasany [kasán] vi cough USAGE May yɔ 'There



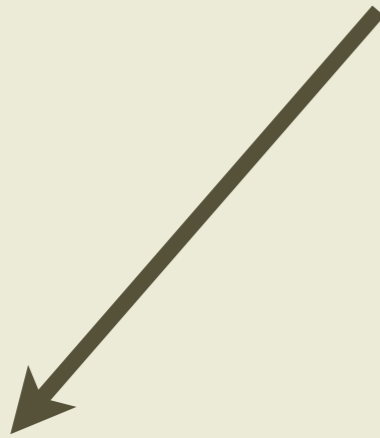
kariz [kari] vt decree good things
kasany [kasap] vi cough USAGE May

Tools

Performance on a very complex source

- FineReader

<i>npf.</i>	<i>pf.</i>	<i>imp.</i>	
བསྐྱེལ་	བསྐྱེལ་	སྐྱེལ་, སྐྱེལ་	
-kyää	-kyää	-kyää, -kyöö	
<i>pres.</i>	<i>fut.</i>	<i>pt.</i>	<i>imp.</i>
[སྐྱེལ་	བསྐྱེལ་	བསྐྱེལ་	སྐྱེལ་]



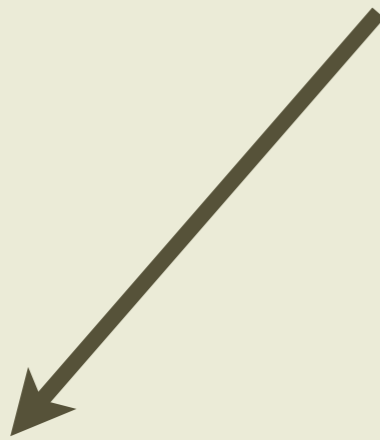
<i>npf.</i>	<i>P/</i>			
"kyaa	"kyaa			
<i>imp.</i>				
"kyaa,	"kyoo			
[gar	qgar	qsp-	gar]

Tools

Performance on a very complex source

- FineReader
("Finnish")

<i>npf.</i>	<i>pf.</i>	<i>imp.</i>	
བསྐྱལ་ -kyää	བསྐྱལ་ -kyää	སྐྱལ་, སྐྱོལ་ -kyää, -kyöö	
<i>pres.</i>	<i>fut.</i>	<i>pt.</i>	<i>imp.</i>
[སྐྱལ་	བསྐྱལ་	བསྐྱལ་	སྐྱོལ་]



<i>npf.</i>	<i>Pf-</i>		
qg^'	"kyää	qspr	"kyää
<i>imp.</i>			
go, -, lq-			
"kyää, "kyöö			
[lar	qgof	qgq-	laj-]

Tools

Performance on a very complex source

- FineReader
(“Finnish”,
HTML
output)

<i>npf.</i>	<i>pf.</i>	<i>imp.</i>	
བསྐྱེལ་ -kyää	བསྐྱེལ་ -kyää	སྐྱེལ་, སྐྱེལ་ -kyää, -kyöö	
<i>pres.</i>	<i>fut.</i>	<i>pt.</i>	<i>imp.</i>
[སྐྱེལ་	བསྐྱེལ་	བསྐྱེལ་	སྐྱེལ་]

<i>npf.</i>	<i>Pf-</i>	<i>imp.</i>
qg^' "kyää	qspr "kyää	go,-, q- "kyää, "kyöö
[ar	qgof	qgq- aj-

Tools

Performance on a very complex source

- Tesseract

<i>npf.</i>	<i>pf.</i>	<i>imp.</i>	
བལྟཱ་	བལྟཱ་	ལྟཱ་, ལྟཱ་	
-kyää	-kyää	-kyää, -kyöö	
<i>pres.</i>	<i>fut.</i>	<i>pt.</i>	<i>imp.</i>
[ལྟཱ་	བལྟཱ་	བལྟཱ་	ལྟཱ་]

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

Tools

Performance on a very complex source

- Tesseract (“Finnish”)

<i>npf.</i>	<i>pf.</i>	<i>imp.</i>	
ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ, ᵛᵛᵛᵛᵛ	
-kyää	-kyää	-kyää, -kyöö	
<i>pres.</i>	<i>fut.</i>	<i>pt.</i>	<i>imp.</i>
[ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ]

npf. pf. imp.
 mgfl' rıgfıı' gfır, šfıı'
 -kyää -kyää -kyää, -kyöö
 pres. fııl. pt. mıp.
 [àfıı' ngfıı' fııgw'
 `f

Tools

Performance on a very complex source

- Tesseract (“Finnish”, HTML output)

<i>npf.</i>	<i>pf.</i>	<i>imp.</i>	
ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ, ᵛᵛᵛᵛᵛ	
-kyää	-kyää	-kyää, -kyöö	
<i>pres.</i>	<i>fut.</i>	<i>pt.</i>	<i>imp.</i>
[ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ	ᵛᵛᵛᵛᵛ]

npf. pf. imp. mgfi' rıgfi' gfir, šfi'

-kyää -kyää -kyää, -kyöö

pres. fi1. pt. mıp.

[àfi' ngfi' fiıgw'

`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