

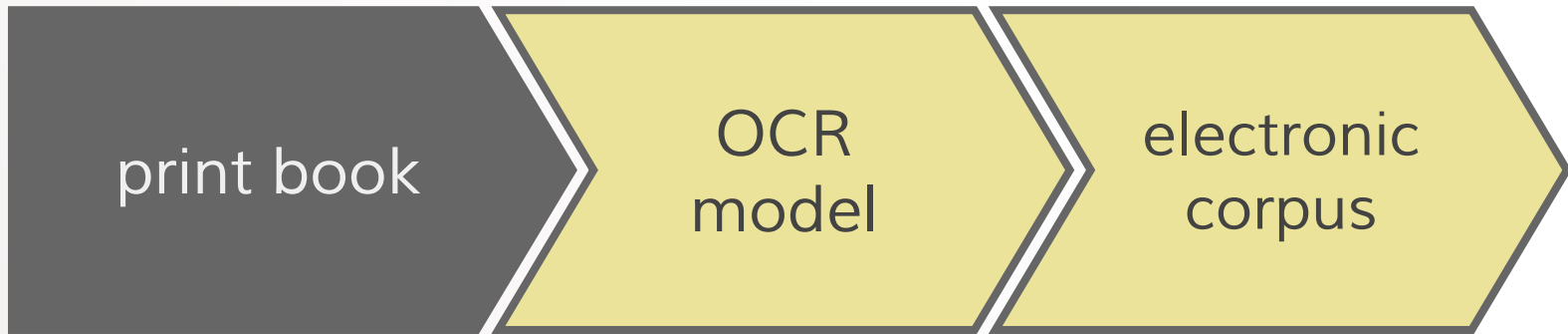
Modern Tech for Old Texts

Training an OCR Model for Northern Haida

Isabell Hubert
with Antti Arppe &
Jordan Lachler

Departmental *Datablitz*
Department of Linguistics
University of Alberta, 18 Sep 2015

The Overall Idea

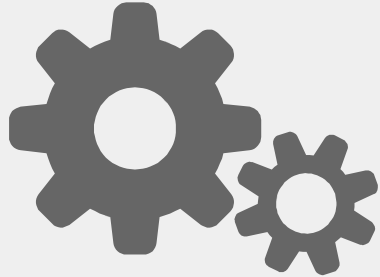


comparing & assessing
different approaches

detailed documentation

OCR System

OCR Engine



- computational work
- language-independent
- *tesseract*

Language Definition/Model

ABC

- “knowledge”
- language-specific
- **this is what we do!**

OCR Model Training Issues

- no model yet
- lots of special characters
- no other OCR model in existence uses the same characters

➔ pick the closest OCR model - in our case: Portuguese

- use this to OCR a couple pages



ã ç ã ç ã ç ã ç

Fix Box Files

- box files define characters computationally
- you teach the model where it “went wrong”

➔ base model

	Char	X	Y	Width	Height
1	l	702	657	12	36
2	'	720	657	6	12
3	s	766	670	17	23
4	ä	786	664	22	29
5	'	809	656	10	15
6	w	821	671	33	22
7	a	856	670	22	23
8	n	881	670	23	22
9	.	909	686	7	7
10	G	975	657	36	36
11	í	1013	659	11	34
12	ñ	1026	661	23	32
13	a	1051	670	22	23
14	'	1075	657	11	15
15	n	1088	670	23	23
16	L	1150	670	22	23
17	!	1176	658	6	36
18	w	1215	671	33	23
19	a	1250	670	23	23
20	'	1274	657	10	15
21	g	1286	671	27	34
22	a	1315	670	22	23
23	n	1340	670	23	23
24	.	1367	687	7	6
25	"	1436	660	14	15
26	H	1461	657	37	35
27	a	1501	670	23	23
28	-	1525	679	13	4
29	i	1539	657	12	35
30	d	704	718	24	36
31	i	730	725	11	28
32	g	770	731	27	35

l' sa'wan.
di gít xēhí'
a'nāñ í'Li:
gí'da dja
Sti'tgA-i l'
Hit!A'n 'a'
LiA gí'sta'a
g' Li gia'

OCR Larger Batches of Pages

- OCR 20-30 pages, validate & fix
- rinse & repeat until you have a whole corpus!

giên la q!ot nañ tc!ā'nan. Gulasqa'ola-i q!ē'idani giên sq!ē'ngua lā'na lā qā'nga-i da gudā'nan. WA'giên Ska'ndal la'ea gi'ā'ldatclan. WA'giên la l' qēngī'si giên da'ñā! A'ñ'ea l' qā'dadjane. T!ā'ngua lā'na ga han isin la l' q!ē-slai'an. La isin l' qā'nan. WA'giên lag^a la l! q!ē'slaiian.

Wa'lū tcāng^a t!a'olē A'ñ'ea l' isda's giên hit!A'n xāg^u l' l!dē'idan. Wa'lū gam la l! kila'āñan ku'na A! l' ti'gai-yani ^an lā l! u'nsatsi A'la. WA'giên l! xao qaođ ha'oisín gulasqa'olē l'ā'ña l! qē'ñidani. Qā'nga-i A! aga'ñ l! xa'ñalane. Gulasqa'oal qē'yiya-i sīñi-djyū'angañan. WA'giên la xēt^g la l! lā l! q!ē'slaiiani la l! qē'ngī'si lū A. WA'giên la l' q!ē'its giên la l' q!ā'da-gaian. "Wā-ā gūs A!ū' l!a'l!a xatca'a-guda'ñwudj. Wēđ han daga'gwañasañ," hñ l' sā'wan.

Silgā'ñ l! stī'ñ!lagalan. Nag^a l!
Swanton (1908), p.184

glen la Y
al-āāi q!ēāā! nañ t ux
āā'ng^aia ani giēā"ā'ñ^a- G
al la',, ^a gudā" sq!ē'ñ ulasqa'
qē^a N ^a gww nan gya ls' o^a
ngl'si . la!dat . WA' , a ĩla lN
T!ā'ñ glēn d' c!an glēnS a
, gua N AñA! . W' . kA'
sra!an laxn^a A'ñ,e,a l, A glēn l n.
la . L ga h qā' a l,
g^a a ĩ A A A da -
Wla l! ,,sl'n l' lsī'n l djane.
- a'lū q'e sLai qa ñan a l, q"
glēn h,, tcān a a.n. . WA' .-e-
lt!A'n g t!ax N glēn
gām la xāg,, l, ole A',,ē
yani ^aA l! kilaw N l!dē'idna l, isd
l! xao n lā L, ěllan ku' an, W' as
l! , qaođ u nsats- na A! l, alū
qē'ñid a'o,x l A'la tī'o .
xa'" ān" lsīn - W o!l-
l. Q.-' A. ' .Ā
djī , e. G añ . qa'ol- glēn
Lū yu An^aa ulasq' gā.l A! e l'ā'ña

giên la q!oi nañ tc!ā'nan. Gulasqa'ola-i q!ē'idani giên sq!ē'ngua lā'na lā qā'nga-i da gudā'nan. WA'giên Ska'ndal la'ea gi'ā'ldatclan. WA'giên la l' qēngī'si giên da'ñā! A'ñ'ea l' qā'dadjane. T!ā'ngua lā'na ga han isin la l' q!ē-slai'an. La isin l' qā'nan. WA'giên lag^a la l! q!ē'slaiian.

Wa'lū tcāng^a t!a'olē A'ñ'ea l' isda's giên hit!A'n xāg^u l' l!dē'idan. Wa'lū gam la l! kila'āñan ku'na A! l' ti'gai-yani ^an lā l! u'nsatsi A'la. WA'giên l! xao qaođ ha'oisín gulasqa'olē l'ā'ña l! qē'ñidani. Qā'nga-i A! aga'ñ l! xa'ñalane. Gulasqa'oal qē'yiya-i sīñi-djyū'angañan. WA'giên la xēt^g la l! lā l! q!ē'slaiiani la l! qē'ngī'si lū A. WA'giên la l' q!ē'its giên la l' q!ā'da-gaian. "Wā-ā gūs A!ū' l!a'l!a xatca'a-guda'ñwudj. Wēđ han daga'gwañasañ," hñ l' sā'wan.

Silgā'ñ l! stī'ñ!lagalan. Nag^a l!

scan

intermediate

final

Assessing Initial Set Size

RQ: How many pages for the initial set are ideal?

- manual labour/
manpower/ accuracy
tradeoff
- comparing models built
from different set sizes to
find sweet spot

84	ê								ê
85	ì								ì
86	í							> í	í
87	î								î
88	ï							<	ï
89	ñ								ñ
90	ô							ü	ô
91	ā							ā	ā
92	ã							ã	ã
93	ē							ē	ē
94	ī								ī
95	ĭ							<	ĭ
96	ı							ı	ı
97	ı̇							ı̇	ı̇
98	ō							ō	ō
99	ū							<	ū
100	ũ							ũ	ũ
101	Ł							Ł	Ł
102	ł							ł	ł
103	.							<	.
104	À							À	À
105	á							á	á
106	â							â	â
107	ã							<	ã
108	ä							ä	ä
109	å							<	å
110	æ							æ	æ
111	“							“	“
112	”							”	”
113	„							<	„
114	Ł							Ł	Ł

Thanks!

References

Swanton, J. R. (1908). Haida Texts - Masset Dialect. In F. Boas (Ed.), *The Jesup North Pacific Expedition, Memoir of the American Museum of Natural History* (Vol. X). Leiden/New York: Brill & Stechert.

tesseract. Open source OCR Engine. Version 3.02.02. github.com/tesseract-ocr

Also, thanks heaps to:

Antti Arppe, Jordan Lachler, Megan Bontogon, Darren Flavelle, Catherine Ford,
Evan Lloyd, Corey Telfer - and John R. Swanton!

template by [SlidesCarnival](#)

photos by [unsplash](#) & [Death to the Stock Photo](#) (license)