

**CEDaR**

# **Developing a Mixed-Methods Pipeline for Community-Oriented Digitization of Kwakwaka Legacy Texts**

Milind Agarwal, Daisy Rosenblum, Antonios Anastasopoulos



# Improving OCR for Kwak'wala Texts

- OCR for legacy data in Kwak'wala
- Why?
  - Rich Written Tradition
  - High value for community research
  - Digital files available only as scanned images

Crabapples.—Wä, la<sup>ε</sup>mē äx<sup>ε</sup>ēdēda ts!Edē<sup>ε</sup>ya, yîxs hē<sup>ε</sup>maē ālēs tselx<sup>u</sup>ts!älaxa tselx<sup>u</sup>wē g'its!âq lāxa lāx·dē gügedzōyosēx la<sup>ε</sup>mē <sup>ε</sup>nāxwaem hē gwēx·<sup>ε</sup>idxa hēlomāg naāgemē lexaxa<sup>ε</sup>ya. Wä, g'il<sup>ε</sup>mēsē lā lōj lāxa tselxwē lā k!adzâlîtaxa lē<sup>ε</sup>wa<sup>ε</sup>yē. k!wāg'alîl lāx hēlk·!ōdenwalîta nānaag lā<sup>ε</sup>wünemas k!wāg'alîl lāx hēlk·!ōdenwalît Wä, laem gēgemxagawalîta laelxa<sup>ε</sup>yē lā dāqē lē<sup>ε</sup>wis lā<sup>ε</sup>wünemē. Wä, lāx·da<sup>ε</sup>xwē lāxa tselxwē qa<sup>ε</sup>s ēp!Exlē māg·înodāla yîsēs hēlk·!ōts!āna<sup>ε</sup>yē. Wä, lā hē dālay

# Project Team



**Milind Agarwal**  
PhD Candidate  
George Mason University



**Antonios Anastasopoulos**  
Assistant Professor (CS, NLP)  
George Mason University



**Daisy Rosenblum**  
Assistant Professor (Indigenous  
Studies, Anthropology)  
University of British Columbia

# Kwakwaka'wakw communities, Kwak'wala language reclamation

- Wakashan language indigenous to North Vancouver Island and opposing mainland in British Columbia (Canada)
- 18 distinct Nations, 5 dialects
- Robust learner community; multiple revitalization efforts.
- >3 orthographies, 2 most common
- Strong tradition of documentation, both internal and external
- Impacted by colonial occupation, forced relocation, removal from land, residential school system



# Kwakwaka'wakw communities, Kwakwala language reclamation

FirstVoices

Dictionary ▾ Learn ▾ Resources ▾ About ▾ Kids

Sign in GUEST

Kwakwala

Search Kwakwala All ▾ 🔍

“”

“Wiga'xan's 'wi'la yaqantala san's yaqandas - Let us all speak our language”

“”

“Dała xan's taxwa'yašan's Kwalskwal'yakw - Carry on with the spirit of our old people”

“”

“Nusansaxw - It belongs to us”

ON THIS SITE

ary Research  
mity Liaison  
anager, GNN

Lands and

BC

GNN School  
Principal, GN

language and  
N School  
age and Cultur

Principal Eke Me

er and Actign

t and Langua

# Value to Language Revitalization

- Manual re-transcription is the norm
- If OCR works well, it can:
  - Reduce manual effort
  - Allow search, increase community access
  - Machine-translation enable NLP applications (language modeling, ASR, transliteration...)



Jaymyn LaVallée, Michayla King  
Gwa'sala-Nak'waxda'xw Language Revitalization Program

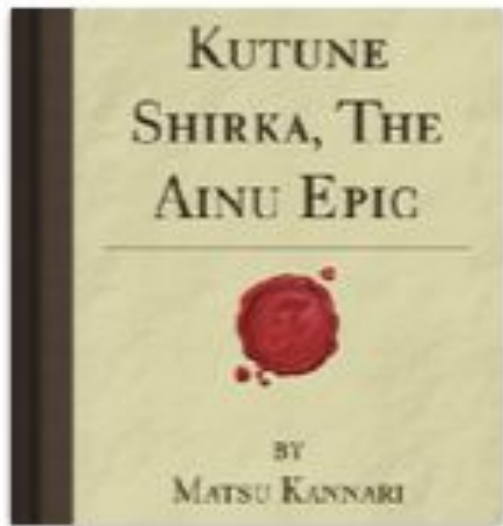


Elder Speakers' Gathering, Gwa'sala-'Nakwaxda'xw Elders' Centre, November 12, 2009

- Allow learners and speakers to spend more time together doing their most important work

# Machine-readable text has many benefits

Matiaxh	Xhunik	jos.om	marim
Mathias	John	work wood (tv).agent	marimb
Matiaxh	Xhunik	wood-worker	(of) n
3.	<sup>14</sup> max.y.al.on	cham —	max
	comp.3.say (tv).PIV	respected man (NC1)	co
	he said		we



*raúusi wa nártúkwa*  
*mí kiltóonaipilikáan nááso*  
*raúusi wa nártúkwa*

Digitizing these texts to machine-readable formats

## Support communities that speak these languages

Make native texts digitally accessible and searchable



Automatic orthography conversion

Aid language researchers, educators, libraries...

## Enable NLP for endangered languages

Annotate datasets for downstream NLP tasks



Build NLP systems to support language revitalization efforts

# Input data

penqwa (e)

pe'nc(a)

pe'nu(a)

peki (a?)

lāxa  
yīsēs

240 BOAS: KWAK'WALA GRAMMAR [TRANS. AMER. PHIL. SOC.]

entrance around neck X 121.38.  
-u'liko in mouth (see -ga into):  
u'liko to throw into mouth III 359.13.  
-u'ndak throat (c = u'ndak): u'ndak speak in throat.  
-ap neck: u'ndak to have on neck III 19.6;  
u'ndak to tuck into neckpiece III 39.3.  
Also: following, behind: u'ndak to stand behind on beach.  
-is u'ndak shoulder and arm above elbow:  
u'ndak to carry on shoulder III 57.16.  
-is u'ndak hand: u'ndak stone-handed III 131.37; u'ndak to put on hand.  
-iso, -ispa chest: u'ndak left side of chest C II 46.17; u'ndak to have hanging on chest N 208.10.  
-ak u'ndak back: u'ndak back blow V 487.4; u'ndak to strike back (max'-)  
Also: u'ndak to follow; u'ndak to roast afterwards; u'ndak to drink afterwards III 41.25.  
-k u'ndak front of body, lap: u'ndak to put in lap V 478.25.  
-u'ndak leg below knee: u'ndak sinews at heel.  
-u'ndak penis: u'ndak with tied---III 138.11 (max'- to tie).  
-is u'ndak shin: u'ndak to put ring around shin III 89.37.  
-k u'ndak knee: u'ndak with scabby knees III 154.11.  
-is u'ndak foot: u'ndak to pinch foot III 96.3; u'ndak to wear shoes CK 281.32.  
-k u'ndak in body, in front of body:  
u'ndak what is in body C II 42.4;  
u'ndak or u'ndak salmon with spawn in body (see -ag).  
-is u'ndak body: u'ndak to sprinkle body III 105.38; u'ndak body III 200.24; u'ndak well grown (tree) V 496.6; u'ndak stone-bodied III 200.9.  
-is u'ndak in front of body: u'ndak to place dish in front of, V 429.23.  
-k u'ndak body (of man, log, etc.) (relating rather to surface of body): u'ndak white-bodied; u'ndak body gets dry all over V 485.6; u'ndak to put on a log III 272.35; u'ndak able-bodied III 208.39; -is u'ndak along a line;  
u'ndak close to a line N 67.96.  
u'ndak to hold (a rope) C 26; 202.97.  
-is u'ndak in mind: u'ndak to feel good III 123.12; u'ndak to think; u'ndak to begin to say in mind (i.e. to think) III 184.3.

4. Limitations of Form  
Generally used with numerals  
u'ndak human beings: u'ndak two persons III 48.21; u'ndak how many persons?  
-u'ndak flat: u'ndak one (day) III 18.2; u'ndak many (leaves) N 298.51.  
-u'ndak long: u'ndak four long ones III 10.12;  
u'ndak number of long ones C III 162.15.  
-u'ndak movement in a long path: u'ndak it goes right through C 26:20.115; u'ndak (to call) only once along street of village C III 218.23.  
-u'ndak round, on surface: u'ndak to hold round thing in mouth C 26:13.6; u'ndak to put down round thing N 485.38; u'ndak woolen blanket (fog on surface) N 691.8.  
-u'ndak tribe: u'ndak five tribes in one village.  
-u'ndak finger-width: u'ndak one finger-width through V 491.6.  
-u'ndak fathom, span: u'ndak seven spans N 110.34.  
-u'ndak times: u'ndak once; u'ndak to stay in house one day.  
-u'ndak days: u'ndak the right number of days III 355.26.  
-u'ndak bundle: u'ndak one bundle N 265.66.  
-u'ndak five pairs of blankets; u'ndak two five pairs (i.e. ten pairs of blankets); also the tenths in each hundred; u'ndak one hundred and ten. Also: u'ndak handsome III 48.29  
-u'ndak pair: u'ndak five pairs (uk'-isema).  
-u'ndak dish: u'ndak two dishes N 516.14; u'ndak one dish V 434.3.  
-u'ndak strings of fish: u'ndak one---  
-u'ndak hole, u'ndak one---; u'ndak size of --- V 332.24.  
-u'ndak layers: u'ndak one---  
-u'ndak one way: u'ndak tide begins to run one way.  
-u'ndak at a time: u'ndak one at a time.  
-u'ndak price (also verb): u'ndak high (many) priced.

5. Temporal Suffixes  
-u'ndak remote past: u'ndak long ago III 12.4;  
u'ndak the late father; u'ndak I came long ago III 142.19.  
-u'ndak recent past: u'ndak plane where he had been III 42.4; u'ndak he went (about a week ago); u'ndak he took a walk.  
-u'ndak future: u'ndak he will have a name III 19.1; u'ndak a future canoe III 85.33.  
-u'ndak transition from present to past:  
u'ndak what he had said III 25.4;

120

al name

slow match

098.20, pengayū

I 49.15,

blow off steam

y steam

Bilingual (English & Kwak'wala)  
Boasian Trilogy: George Hunt & Franz Boas collaboration  
1897-1948  
over 11 published volumes (~10000 pages)

Additional trove of unpublished archival text (type- & hand-written)



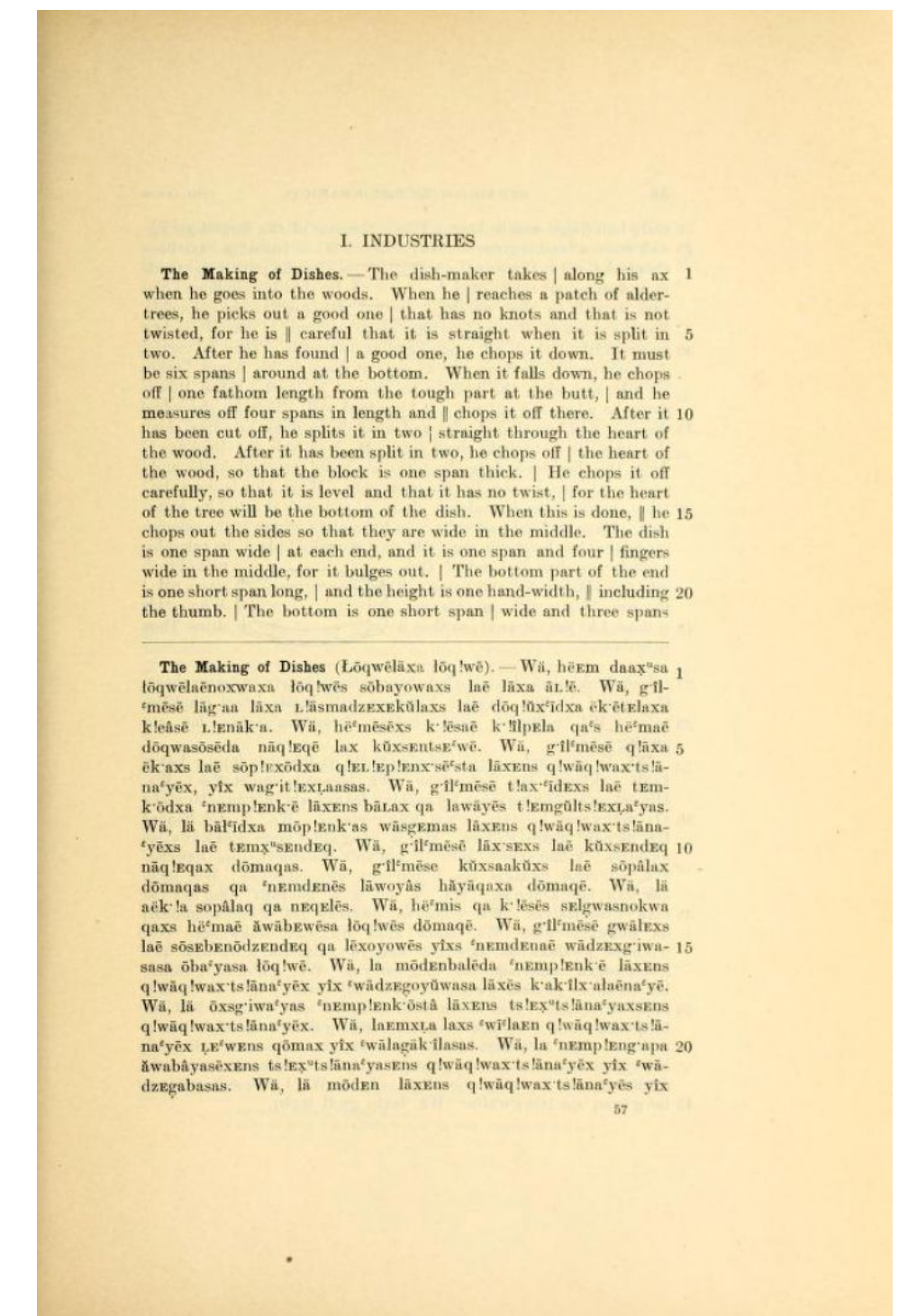
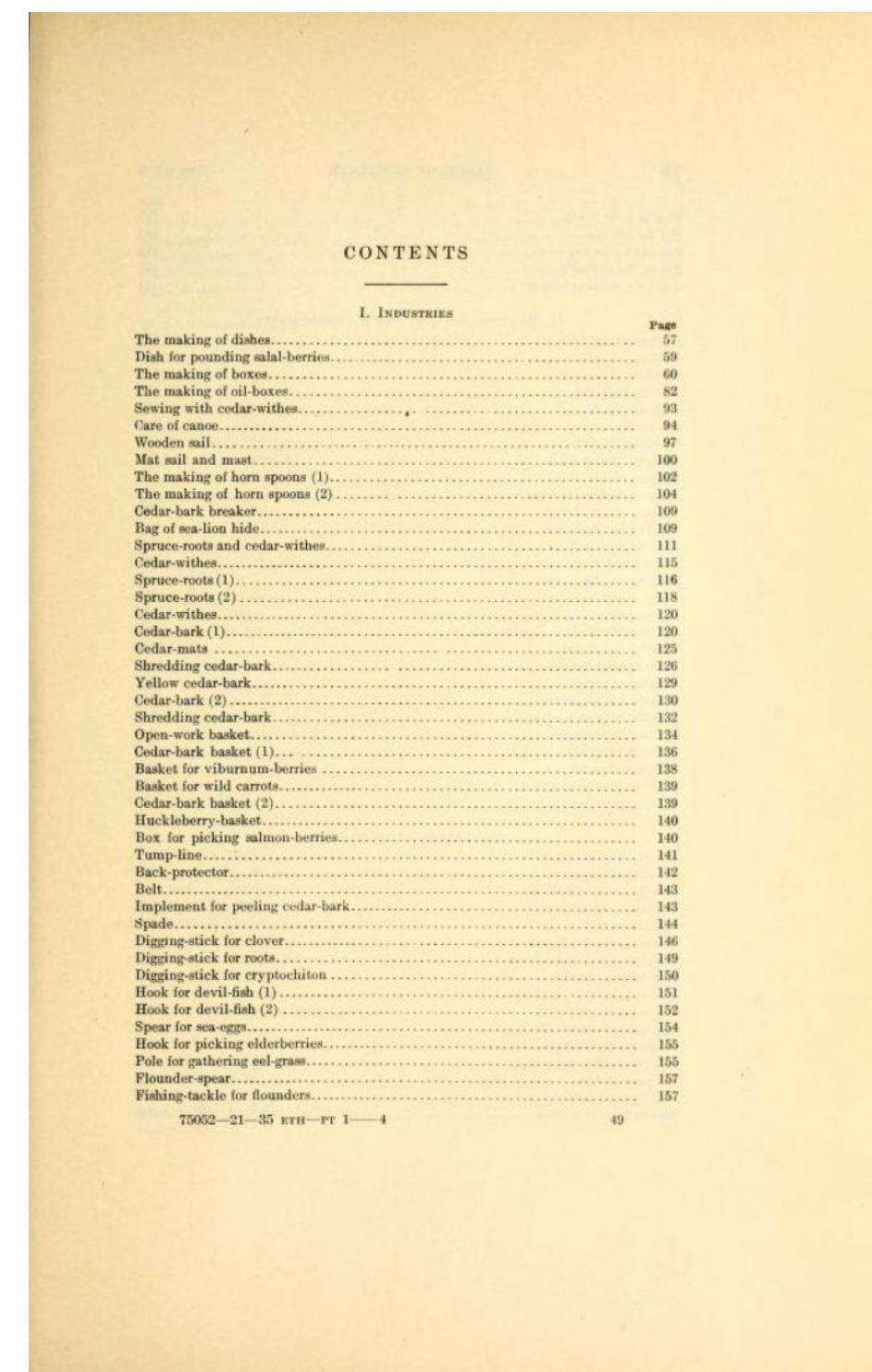
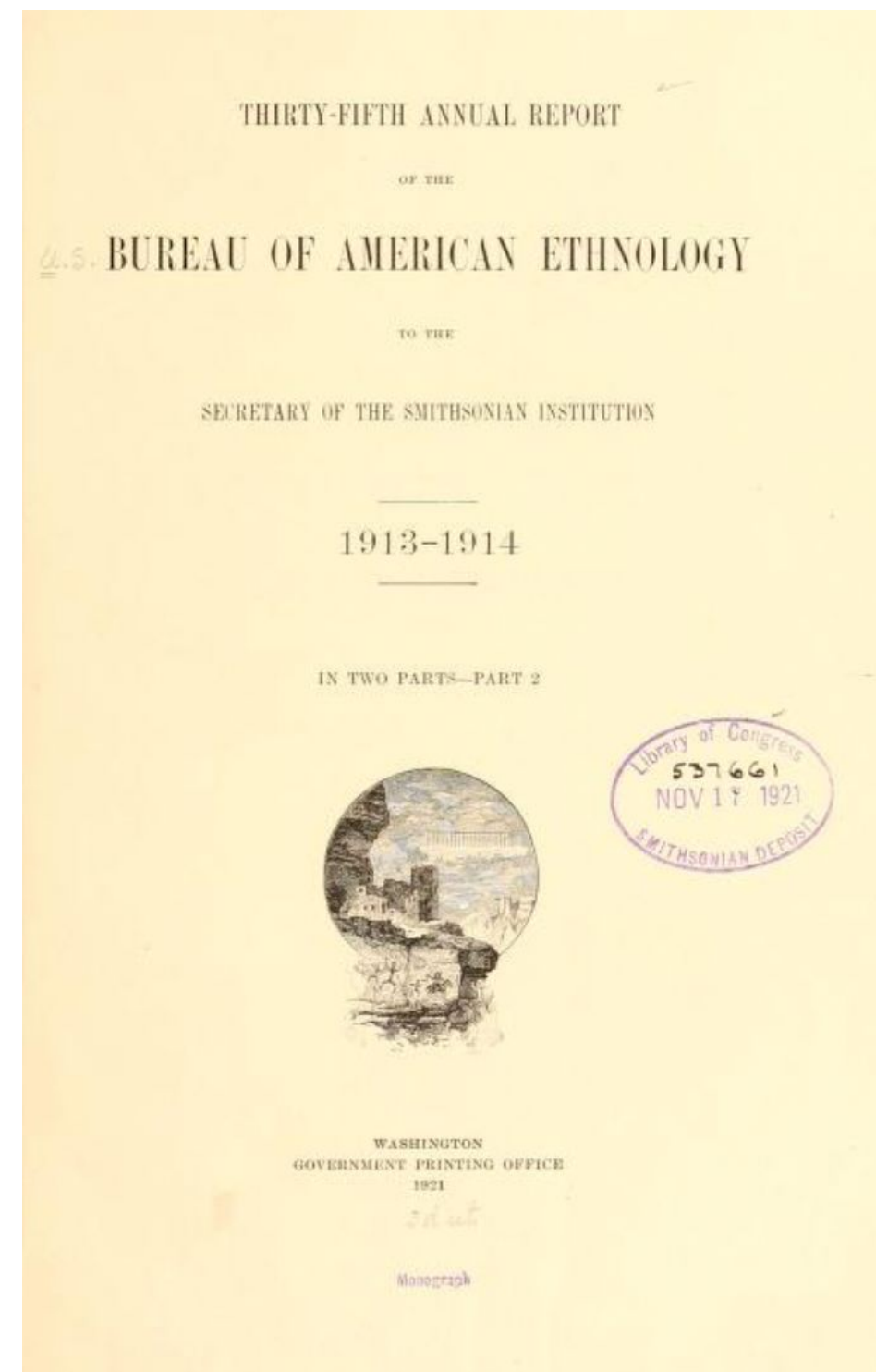
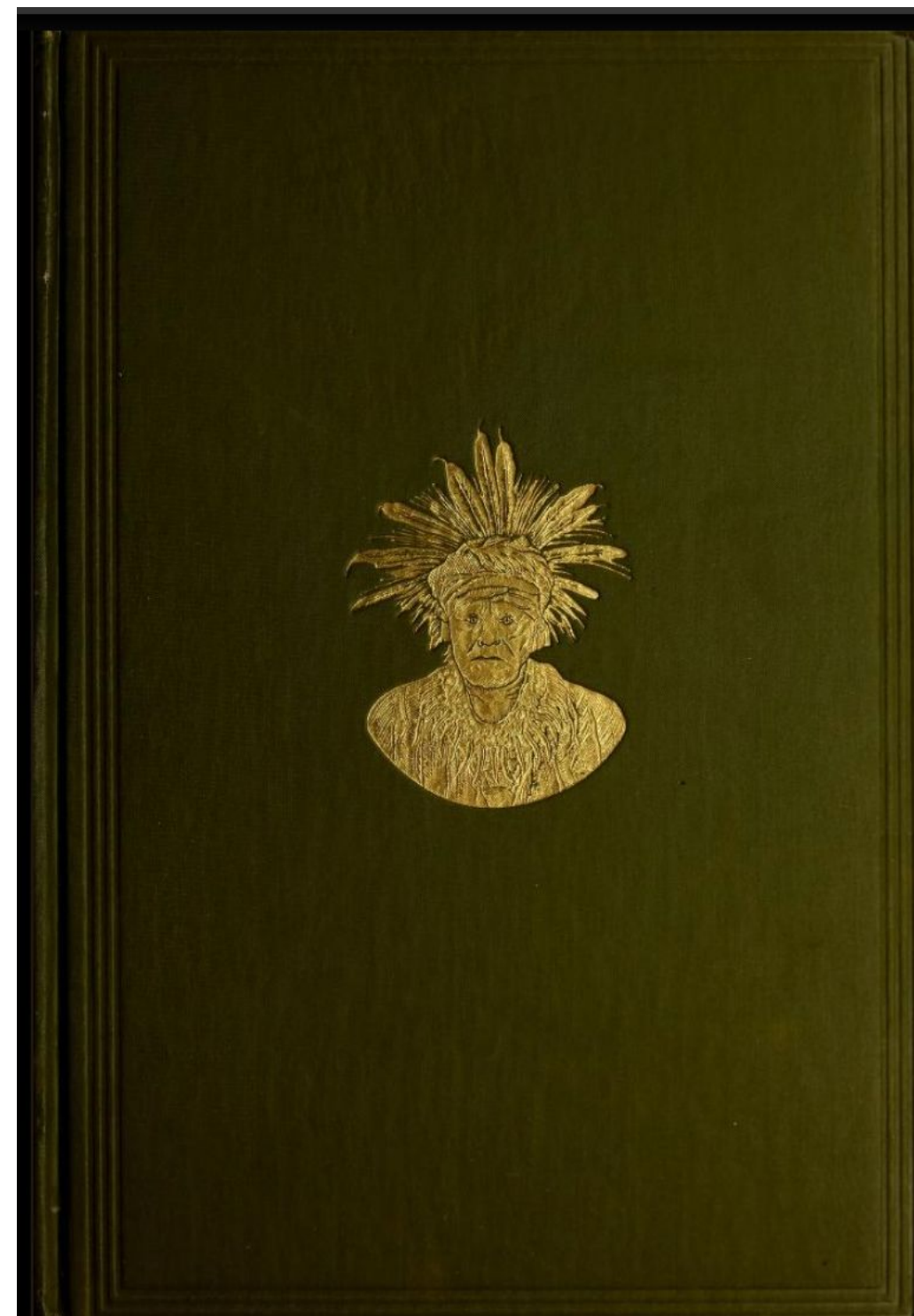
# Initial Challenges

la<sup>ε</sup>mē<sup>ε</sup> nāxwāem hē gwēx<sup>ε</sup>īc  
naāgēmē lEXa<sup>ε</sup>ya. Wä, g·îl<sup>ε</sup>  
lāxa tselxwē lā k·!adzâlîlaxa  
k!wāgalîl lāx hētk·!ōdenwalîl

- Existing OCR systems don't work very well on many under-resourced languages
- Even pretrained OCR systems that support 30+ scripts do not support the early 20th c ethnographic writing systems
- More difficult on old books and low-quality scans
- Difficult to train a good model without data!

# Text Selection: Community-determined

Boas, Franz, and George Hunt. 1921. "Ethnology of the Kwakiutl. Based on Data Collected by George Hunt." In *Thirty-Fifth Annual Report of the Bureau of American Ethnology Presented to the Secretary of the Smithsonian Institution*, parts 1 and 2, 42–794, i–xi [index]. Washington, DC : Government Printing Office.



# First Phase

- Post-correction OCR model
- CER improved from 30% > 4%!
- Automated conversion from Boas > U'mista
- 1401 separate txt files (66.txt, 67.txt...)
- Compiled first draft PDF for distribution

## I. INDUSTRIES

The Making of Dishes. The dish-maker takes | along his ax when he goes into the woods. When he | reaches a patch of alder- trees, he picks out a good one that has no knots and that is not twisted, for he is || careful that it is straight when it is split in two. After he has found a good one, he chops it down. It must be six spans around at the bottom. When it falls down, he chops off one fathom length from the tough part at the butt, | and he measures off four spans in length and || chops it off there. After it has been cut off, he splits it in two | straight through the heart of the wood. After it has been split in two, he chops off the heart of the wood, so that the block is one span thick. He chops it off carefully, so that it is level and that it has no twist, | for the heart of the tree will be the bottom of the dish. When this is done, || he chops out the sides so that they are wide in the middle. The dish is one span wide | at each end, and it is one span and four | fingers wide in the middle, for it bulges out. | The bottom part of the end is one short span long, | and the height is one hand-width, || including the thumb. The bottom is one short span wide and three spans

The Mak-ing oo Dihg.). Wä, hēm daax<sup>sa</sup> lōqwēlaēnoēwaxa lōq!wēs sōbayowaxs laē lāxa ālē. Wä, g-īl- <sup>mēsē</sup> lāg-aa lāxa lāsmadzexēkūlaxs laē dōqūx-<sup>īdxa</sup> ēk-ētēlaxa k!eāsē l!ēnāk-a. Wä, hē<sup>mēsēxs</sup> k!ēsāē k!īpēla qa<sup>s</sup> hē<sup>maē</sup> dōqwasōsēda nāq!ēqē lāx kūsxsntsē<sup>wē</sup>. Wä, g-īl<sup>mēsē</sup> q!āxa ēk-axs laē sōp!ēxōdxa q!ēl!ēp!ēnx-se<sup>sta</sup> lāxens q!wāq!wax-ts!ā- na<sup>yēx</sup>, yix wāg-it!ēx,aa<sup>sas</sup>. Wä, g-īl<sup>mēsē</sup> t!ax-<sup>īdēxs</sup> laē tēm- k!ōdxa <sup>nēmp!ēnk-ē</sup> lāxens bālxax qa lawāyēs t!ēmgūlts!ēx,ā<sup>yas</sup>. Wä, lā bāl<sup>īdxa</sup> mōp!ēnk-as <sup>wāsgēmas</sup> lāxens q!wāq!wax-ts!āna- <sup>yēxs</sup> laē tēm<sup>sendēq</sup>. Wä, g-īl<sup>mēsē</sup> lax-sexs laē kūsxsndēq nāq!ēqax-dōmaqas. Wä, g-īl<sup>mēsē</sup> k!ūxsakūxs laē sōpālx dō- maqas qa <sup>nēmdēnēs</sup> lāwoyās hayaqaxa dōmaqē. Wä, lā aēk!a sōpālaq qa nēqēlēs. Wä, hē<sup>mis</sup> qa k!ēsēs selgwasnōkwa qaxs hē<sup>maē</sup> āwābē<sup>wēsa</sup> lōq!wēs dōmaqē. Wä, g-īl<sup>mēsē</sup> gwālēxs laē sōsebenōdzēndēq qa lēxoyowēs yixs <sup>nēmdēnāē</sup> <sup>wādze<sup>xg</sup>-iwa-</sup> sasa ōba<sup>yasa</sup> lōq!wē. Wä, la mōdēn- balēda <sup>nēmp!ēnk-ē</sup> lāxens q!wāq!wax-ts!āna<sup>yēx</sup> yix <sup>wādze<sup>goyu</sup>wasa</sup> lāxēs k-āk-īlxalāēna<sup>yē</sup>. Wä, lā ōxsg-iwa<sup>yas</sup> <sup>nēmp!ēnk-ōstā</sup> lāxens ts!ēx<sup>ts!</sup>āna<sup>yaxsens</sup> q!wāq!wax-ts!āna<sup>yēx</sup>. Wä, laēm<sup>x,ā</sup> lāxs <sup>wā</sup> laēm q!wāq!wax-ts!ā- na<sup>yēx</sup> tē<sup>wins</sup> qōmax yix <sup>wāla<sup>gak</sup>-īlasas</sup>. Wä, la <sup>nēmp!ēng-apa</sup> āwābā<sup>yasēxens</sup> ts!ēx<sup>ts!</sup>āna<sup>yasens</sup> q!wāq!wax-ts!āna<sup>yēx</sup> yix <sup>wā-</sup> dzēgababas. Wä, lā mōdēn lāxens q!wāq!wax-ts!āna<sup>yēs</sup> yix

The Making o'o Dihg.). We, he'am da'axwsa luḡwila'ino'iwaḡa luḡ'wis subayuwaḡs la'i laḡa ali. We, gil' misi laga'a laḡa tlesmadzaxakulaḡs la'i duḡux'īdxa ekitaḡaḡa ke'osi t'ḡanaka. We, he'misiḡs kisa'i kīpāla ḡa's he'ma'i duḡwasusida naḡḡāki laḡ kuḡsantsa'wi. We, gil'misi ḡaxa ekaḡs la'i supaxudḡa ḡat'ḡapānxsā'sta laḡans ḡ'waḡ'waxt'sa- na'yix, yix wagitaxdīa'asas. We, gil'misi tax'īdaxs la'i tām- kudḡa <sup>nāmpānki</sup> laḡans batlax ḡa laweyis tām<sup>gult</sup>'saḡdīa'yas. We, le bal'īdxa muḡānkas <sup>wasḡa-</sup> mas laḡans ḡ'waḡ'waxt'sana- <sup>yixs</sup> la'i tām<sup>xwsā</sup>ndaḡ. We, gil'misi lāxsaxs la'i kuḡsāndaḡ naḡḡāx- dumaḡas. We, gil'misi k!uḡsa'akuḡs la'i supolaḡ dumaḡas ḡa <sup>nāmdānis</sup> lawuyos hayaḡaxa dumaḡi. We, le a'ēka supolaḡ ḡa naḡḡālis. We, he'mis ḡa kisis saḡgwasnukwa ḡaxs he'ma'i ḡwobā'wisa luḡ'wis dumaḡi. We, gil'misi ḡwalaḡs la'i susāḡānudzāndaḡ ḡa liḡuyuwis yixs <sup>nāmdāna</sup>'i <sup>wadzaḡgiwa-</sup> sasa uba'yasa luḡ'wi. We, la mūdānbalida <sup>nāmpānki</sup> laḡans ḡ'waḡ'waxt'sana'yix yix <sup>wadzaḡgiyuwasa</sup> laḡis kakilḡala'ina'yi. We, le uḡsḡiwa'yas <sup>nāmpānkusto</sup> laḡans t'saxwt'sana'yaxsāns ḡ'waḡ'wax- t'sana'yix. We, la ḡmḡdīa laḡs <sup>wa</sup> la'ām ḡ'waḡ'waxt'sa- na'yix dīa'wīns ḡumaḡ yix <sup>walaḡakīlasas</sup>. We, la <sup>nāmpāngapa</sup> ḡwobā'yasiḡans t'saxwt'sana'yaxsāns ḡ'waḡ'waxt'sana'yix yix <sup>wa-</sup> dzēgababas. We, le mudān laḡans ḡ'waḡ'waxt'sana'yis yix

# Communities and Consultations

First draft of text distributed widely to ~50 community experts and areal knowledge specialists/researchers for feedback

- GN Language Program
- Staff of the K̓wala adult immersion language program
- Language teachers in the GNN Elementary School
- Heritage language learners and teachers at UBC
- Elders working with the language
- Academic researchers, linguists

Draft OCR Transcription of Ethnology of the Kwakiutl (Boas & Hunt 1921, Smithsonian) Survey <https://docs.google.com/forms/u/1/d/15CG9Nraz-OjJ1e499ieEiu4eVmh...>

## Draft OCR Transcription of *Ethnology of the Kwakiutl* (Boas & Hunt 1921, Smithsonian) Survey

\* Indicates required question

1. Do you like seeing both orthographies, Boas-Hunt and U'mista, on the same page? \*

Mark only one oval.

- Yes  
 No

2. Of the options listed below, which would you prefer? (Please select all that apply) \*

Check all that apply.

- Just U'mista and English  
 Just Hunt-Boas and English  
 U'mista, Hunt-Boas, and English  
 Other: \_\_\_\_\_

3. What errors, and types of errors, do you notice in the U'mista writing? Which ones are the most distracting for you? \*

---

---

---

---

---

# Secondary Challenges

## • Revising 1921

- Restore line numbers
- Restore missing titles

How to make this model transferable to other texts?

Essential:

- How to handle diverse formatting, layouts & arrangement of English and Kwak'wala?

Nice to have:

- Images & Figures

20.99; *ō'gwīwūlē* (*ō-g'iu-lē*) forehead on water; bay in front of *Xūmdasbē*, Hope Island — M 370, 391; *ō'gwitēmē* (*ō-g'it-[g]ēm-ē*) head of body of round thin 22.135; *ō'gwitēmēs nō'mas* (*ō-g'it-g' body of Old Man 7.112; ō'gwitēmēs gag'ēwas* 6.30; *ō'gwitēmālis* (*ō-g'it-g' of body 8.15; 8a.70; ō'gwitēmēlē* (*ō-g' on water 10.59; ōk'ū' nē* (*ō-k'ūm-ē*), Deserter's Island 6.47 — CXXVI 220; beach body, sand bar in river 13.50 CXXVI 104; *ō'gūmālis* (*ō-gēm-a-l-is*) (*ō-gēm-l-!a*) front rock 3.47 — R 1221. *is*) beach at river mouth 7.59; *ōx'sī* river 8.109; 22.136; *ō'x'sīdzē* (*ō-x'sīs= 22.67; ō'x'sīdza'laa* (*ō-x'sīs-a-l-!a*) rock (*ō-x'ts!ē(?) -l-!a*) front of long, steep *ō'x'sā* (*ō-x'sā*) passage through 11.32; *ōx'sā* through 10.74; 14.4; *ō'x'sālēsēla* (*ō-x's through 14.115; ōxā'wē* (*ō-xo-ē*) nec — III 149.22; *ōxsdē* (*ō-xsd-ē*) hind 14.9; 18.103; 20.27; *ōxsde'las* (*ō-xsd- inside CXXVI 95; ōxsde'lis* (*ō-xsd- 8.3, 18; 21.4; 22.134; ō'xsdēlē* (*ō-x 10.119; ōxsdēlēs nō'mas* (*ō-xsd-e-lē-s of Old Man 7.108; ō'x'ulā* top of head 2: head beach 3.107; 18.7; 20.43 — (*ō-x'la-ato-a-l-is*) beach at hind end 20.89; *ōx'la'laa* (*ō-x'la-l-!a*) rock : 16.50; 17, between 15 and 16; *ōx'la'* hind end 3.6; 6.93; 8.10, 19; 10.28, *ō'yaa* (*ōya-a*) (*ōya Nak'wax'da'x'*; *gw rock 7.51; 10.69; ō'yāēlē* (*ōya-ēl-b 10.215; ō'yabaa'* (*ōya-b-!a*) outside r (*ōya-b-ē*) outward point 18.9; *ō'yagē* outward rock 4.7; 15.61; 17, between *ō'ya-xsta'his* (*ōya-āxst-a-l-is*) beach : *ō'bēk'* (*ōp-ku?*) whispered (?) 15.143 (*ō'manis*), Nootka name; 4.63 — M 329 *ō'maxsdelis* (*ōm-exsd-el-is*) 13.36, 43 *ōmaxsēla* 13.25 *ō'sēq'* (D 104 *osē-ka*) said to mean “grey trees 11.6 — M 677.17; X 229.18; 1388.62; *āō'saaqūm* (D 105 *ā-ows-a- 14.128; ō'saaqūmlis* 16.3

TTL. 473  
 5'ts lā'xa grī'dasē yix la grītsle-

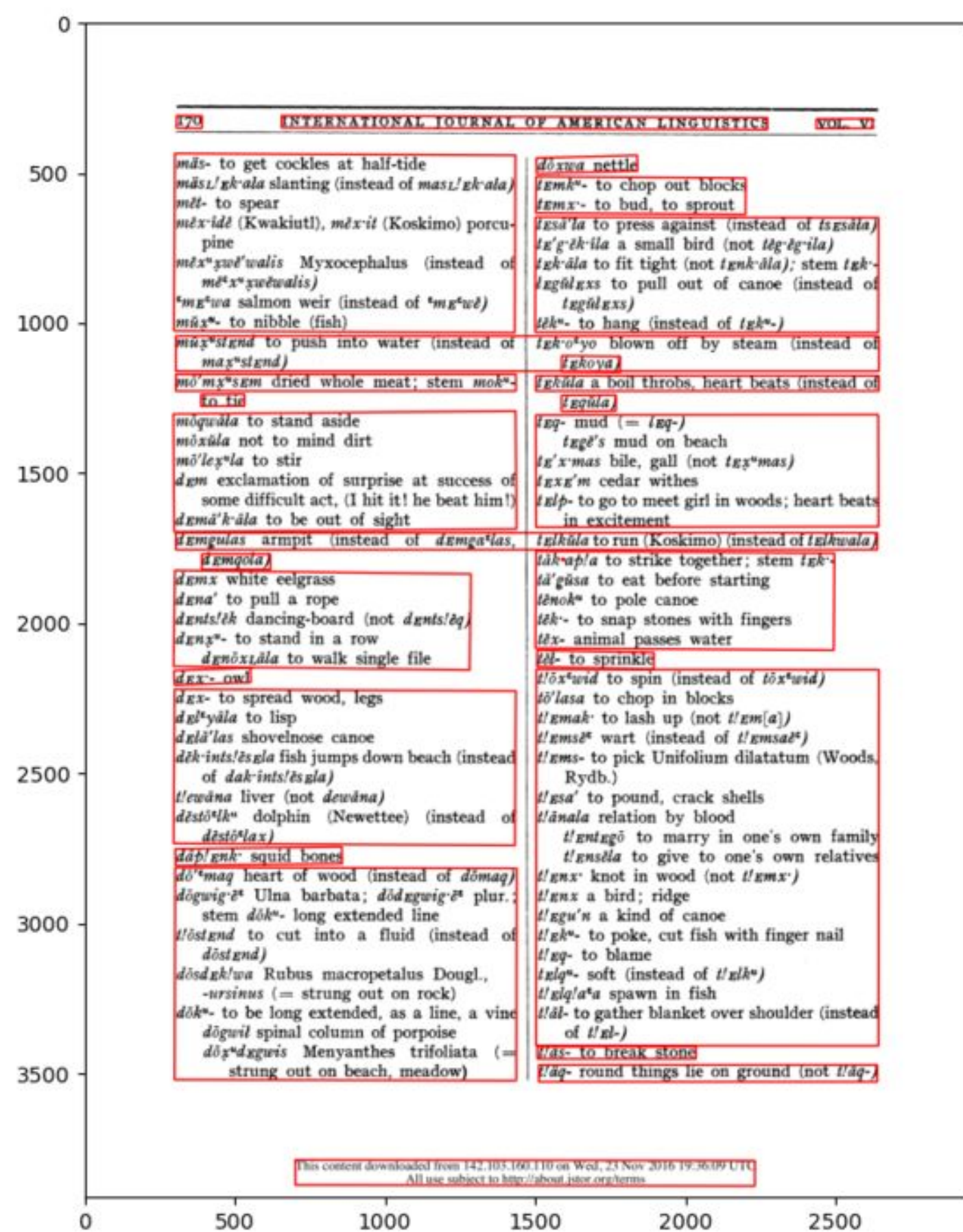
the speaker of the house said, “Come here, Mouse, and go again to see. Now go really, and try to see what affects us, and what is the reason that this death-bringer does not act in the right way.” Then the Mouse went out of the winter-dance house of the mountain-goats. She went at once to the hiding-place of TEWĪ'x'ī'lak<sup>1</sup>, and said, “O friend! take care. When they sing again, you must enter. They will name you at once Dā'bend.<sup>1</sup> When you jump in at the door of the house, you must take hold of that feather, and say while you are holding it, ‘I am Dā'bend.’ Then the large man will let go of it, but you must hold the feather. That is (what I want to say).” Then the Mouse went back into the house, and said, “O friends! I have been all round our world.” Thus said the Mouse.

á'Emxat! qlwē'fēda. Wā, lā'laē ē't!lēd yā'q!eg'atēda yā'yaq!ente'mē'fasēda g'ō'kwē. 1  
 only again they stopped. Well, then it is again spoke the speaker of the house of the house.  
 Wā, lā'laē nē'k'a : “Gē'lag'a, Hā'la'mālag, qa's lā'ōs ē't!lēd dō'x'wida. Wē'gra 5  
 Well, then it is he said : “Come here, Mouse (woman), that you go again to see. Go on  
 said  
 á'lax'īd dō'q'wax hā'tā grā'xens; lā'grī'asik' k'lēs hē'fē'lag'a da hālā'yuk.”  
 really look for what affects us; the reason for not being right this death-bringer this.  
 Wā, hē'x'īdaem'la'wisē Hā'la'mā'laga la qa'sēd qa's lā lā'wels lā'xa  
 Well, at once it is said the Mouse (woman) then walked, that she went to go out at the  
 ts'lā'gats'lā'sēda mē'mē'lxlowē. Wā, lā'laē hē'nā'kulaem lāx wu'ndzasas 5  
 ts'ē'ts!ēqa house of mountain-goats. Well, then it is she went at once to the hiding-place  
 said  
 TEWĪ'x'ī'lakwē. Wā, lā'laē nē'k'a : “ēya, qāst, wē'gra yā'l'lāx, lā'ems lāl  
 TEWĪ'x'ī'lak'. Well, then it is she said : “O friend! go on take care, then you will  
 said  
 lā'ē'lōl qa'xō ē't!lēl dē'nx'ēlō. Hē'mā'qō lāl lē'x'ēlxēs lē'gemōsē Dā'bendē,  
 you will when will again will sing. At once they will name your your name Dā'bend,  
 enter, will will will  
 wā, lā'LES dēwī'l lā'xwa tlēx'ī'lāxsa g'ō'kwēx. Wā, lā'LES dā'x'īlxwa  
 well, then you jump in at the door of the house. Well, you will will take the  
 will  
 ts!ē'lts!ēlk'ēx; lā'LES nē'x'lōl : “Nō'gwaem Dā'bend,' qa'sō lāl dā'laleq.  
 feather here; then you you will you will ‘I am Dā'bend,' when then hold it.  
 will say : will will will  
 Wā, hē'x'īdaem'la'wisōxda wā'lasēx begwā'nem mēx'ē'LEQ". Wā, lā'LES 10  
 Well, at once will this large this man will let go of it. Well, then however  
 you will  
 dā'lax'sāem'la'xa ts!ē'lts!ēlk'ē. Hē'mēq." Wā, lā'laē qa'sēdē Hā'la'mā'laga  
 will but only hold the feather. That is it." Well, then it is walked the Mouse (woman)  
 said  
 qa's lē lā'ē'l lā'xa g'ō'kwē. Wā, lā'laē nē'k'a : “ēya nē'nemōkwā'ī, lā'mx'dēn  
 that went to in the house. Well, then it is she said : “O friends! I have  
 she enter said  
 lē'ēstāl'sēla lā'xwa awī'ēstāxsens nā'lax," nē'x'laē Hā'la'mā'laga.  
 gone around the at this around of our world," said it is said the Mouse (woman),  
 world

<sup>1</sup> That means “to take hold of end.”



# First-Pass OCR



## Google Cloud Vision OCR - Text Detection Module

- Advantages
  - Free to use, up to 1000 pages per month
  - Access through Google Cloud API services
  - Quick and easy setup, and rich metadata where available
- Disadvantage:
  - Does not allow finetuning
  - Language ID labels are often missing
  - Bounding boxes may not be structurally coherent

# LangID

Stretched rom therow right through laid oxwalelodayu lax ogiwa la  
To the stern but when the little he benda la lax o xlayas lala -  
Canoe is upright it is this way )) qexs heneda xwaxwagle  
ga gwa lega ).



English-language detector (fastText)

Stretched rom therow right through laid oxwalelodayu lax ogiwa la  
To the stern but when the little he benda la lax o xlayas lala -  
Canoe is upright it is this way )) qexs heneda xwaxwagle  
ga gwa lega ).



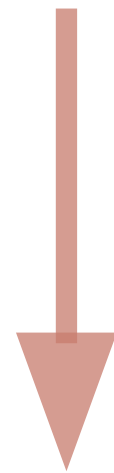
# LangID

fastText custom language identification training

- bilingual first-pass texts from Boas collection used for training
- 1000 sentences for each language (Kwak'wala and English)
- Trained using default fastText parameterizations
- Achieves 99+% accuracy on held-out test set
- We apply this to the first-pass text to potentially use language as a proxy for structure

# Masking

Stretched rom therow right through laid oxwalelodayu lax ogiwa la  
To the stern but when the little he benda la lax o xlayas lala -  
Canoe is upright it is this way )) qexs heneda xwaxwagle  
ga gwa lega ).



High-likelihood English spans are temporarily masked

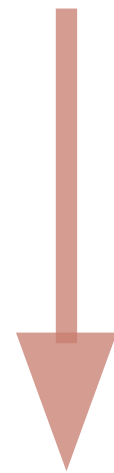
laid oxwalelodayu lax ogiwa la  
benda la lax o xlayas lala -  
)) qexs heneda xwaxwagle  
ga gwa lega ).

# Masking

- Texts are diversely formatted with illustrations, figures, page numbers, line numbers etc.
- The post-correction model requires only Kwak'wala text and doesn't handle non-Kwak'wala text very well
- We apply a masking layer that temporarily hides:
  - line numbers (that provide cross-reference information to other collections (very important for readers))
  - page numbers
  - punctuation
- Done with a simple deterministic Python script

# Post-Correction

laid oxwalelodayu lax ogiwa la  
benda la lax o xlayas lala -  
) qexs henaleda xwaxwagle  
ga gwa lega ).



Remaining Kwak'wala text is sent for post-correction,  
as per Rijhwani et al. (2020)

la'e do'x<sup>3</sup>waLElōdayu lāx o'g'iwa<sup>3</sup>yas la  
hē'bendāla lax ō'xḷa<sup>3</sup>yas, lāḷa-  
qēxs hē'naḷaēda xwā'xwagumLē  
ga gwälega).

# Post-Correction

- Can automatically correct errors in very low-resource OCR settings
- Train a very small correction model on a sample of first-pass and gold reference sentences
- Multi-source neural architecture (based on Rijhwani et al. 2020) which was shown to reduce OCR character-level errors by 30-60%
- We use the model as is, with the lexically-aware decoding setting turned off as it was shown not to benefit Kwak'wala
- We train the model from scratch, replicate the CER results from the original paper, and then apply the model to our dataset

# Reconstruction

la'e do'x<sup>3</sup>waLElōdayu lāx o'g'iwa<sup>3</sup>yas la  
hē'bendāla lax ō'x<sub>7</sub>ḷa<sup>3</sup>yas, lāḷa-  
qēxs hē'nāḷaēda xwā'xwagumLē  
ga gwälega).

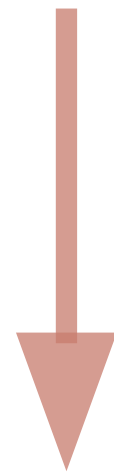


English text and other masked tokens are reintroduced at the appropriate indices. Orthographies can be converted (U'mista or SD-72)

Stretched rom therow right through  
To the stern but when the little  
Canoe is upright it is this way  
la'e do'x<sup>3</sup>waLElōdayu lāx o'g'iwa<sup>3</sup>yas la  
hē'bendāla lax ō'x<sub>7</sub>ḷa<sup>3</sup>yas, lāḷa-  
qēxs hē'nāḷaēda xwā'xwagumLē  
ga gwälega).

# Evaluation

la'e do'x<sup>3</sup>waLElōdayu lāx o'g'iwa<sup>3</sup>yas la  
hë'bendāla lax ō'x<sub>ɾ</sub>a<sup>3</sup>yas, lā<sub>ɾ</sub>a-  
qēxs hë'na<sub>ɾ</sub>aēda xwā'xwagumLē  
ga gwälega).



Compare final corrected outputs to original gold labels, at the character level to obtain CER

la'ē dō'x<sup>s</sup>waLElōdayu lāx ō'g'iwa<sup>s</sup>yas la  
hë'bENDāla lāx ō'x<sub>ɾ</sub>a<sup>s</sup>yas, lā'<sub>ɾ</sub>a-  
qēxs hë'na<sub>ɾ</sub>aēda xwā'xwagumLē  
g'a gwä'lēg'a).





# Results

- Community-oriented prioritization of texts for digitization
- With our mixed-methods pipeline, including language identification, masking, and automatic post-correction, we achieve for our Kwak'wala texts:
  - ~50% decrease in character error rate
  - ~87.5% reduction in structural error
    - insertion, deletion, and maximal move operations required across the output page to make it identical with the reference text (Kanai et al, 1995)

	Jesup 5.1, 1902		Kwakiutl, 1909	
	CER	SER	CER	SER
<b>First Pass</b>	0.43	25	0.33	18
<b>Corrected</b>	0.18	2	0.15	3

Table 1: For both books, we find that using our pipeline greatly reduces not only textual errors (CER) but also greatly improves the layout and structure (SER)

# Next Steps

- Stable OCR Packages for Kwak'wala:
  1. Developer friendly (repeatable & transferable)
  2. End-user friendly
- Investigation of multiple langID modeling approaches for Kwak'wala (other than fastText)
- Create more gold reference texts for better evaluation
- More extensive benchmarking across more books/collections
- With consent of the language community, share digitized texts with the data hosting institutions, such as the American Philosophy Society and Columbia University Rare Books and Special Collections

# Gilakas'la! Thank you!

Developing a Mixed-Methods Pipeline for Community-Oriented Digitization of Kwak'wala Legacy Texts

Milind Agarwal, Daisy Rosenblum, Antonios Anastasopoulos



GitHub Repository  
Contains code for all stages of the pipeline, as applied on Kwak'wala



# Next Steps

- Text extraction from two volumes (Boas and Hunt, 1921)
- Transliteration into two community-preferred orthographies (U'mista, Liq'wala)
- Apply model to other languages written in same typeface
- Develop/improve OCR model for unpublished typescript
- Create dataset for developing other NLP tools