

Neural Polysynthetic Language Modelling





2019 | INTERNATIONAL YEAR OF
Indigenous Languages



Background



Intersecting machine learning & linguistic fieldwork

St. Lawrence Island Yupik



Intersecting machine learning & linguistic fieldwork



Moving forward



Lane Schwartz



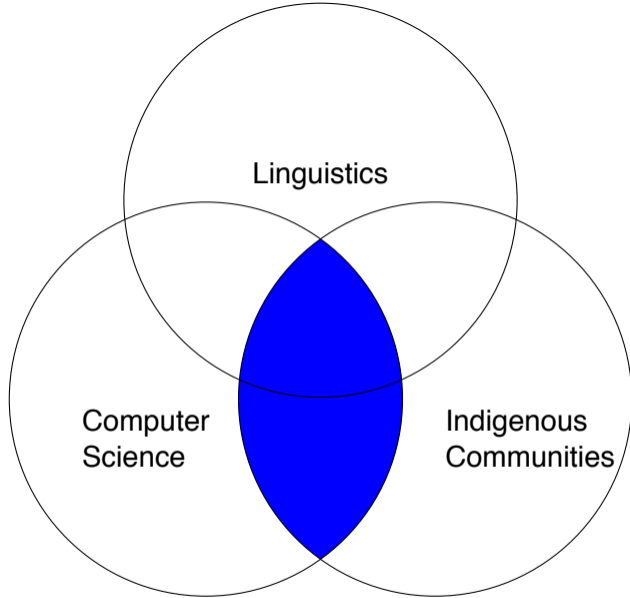
2019 | INTERNATIONAL YEAR OF
Indigenous Languages

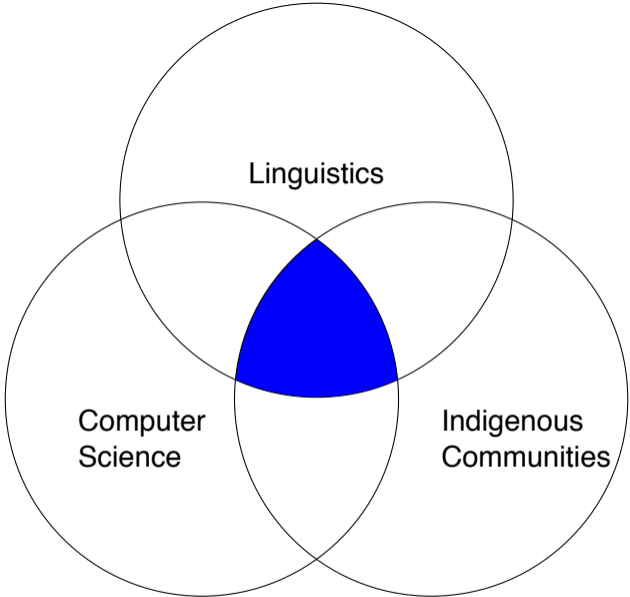
- Increasing understanding, reconciliation and international cooperation.
- Creation of favorable conditions for knowledge-sharing & dissemination of good practices with regards to indigenous languages.
- Integration of indigenous languages into standard setting.
- Empowerment through capacity building.
- Growth and development through elaboration of new knowledge.

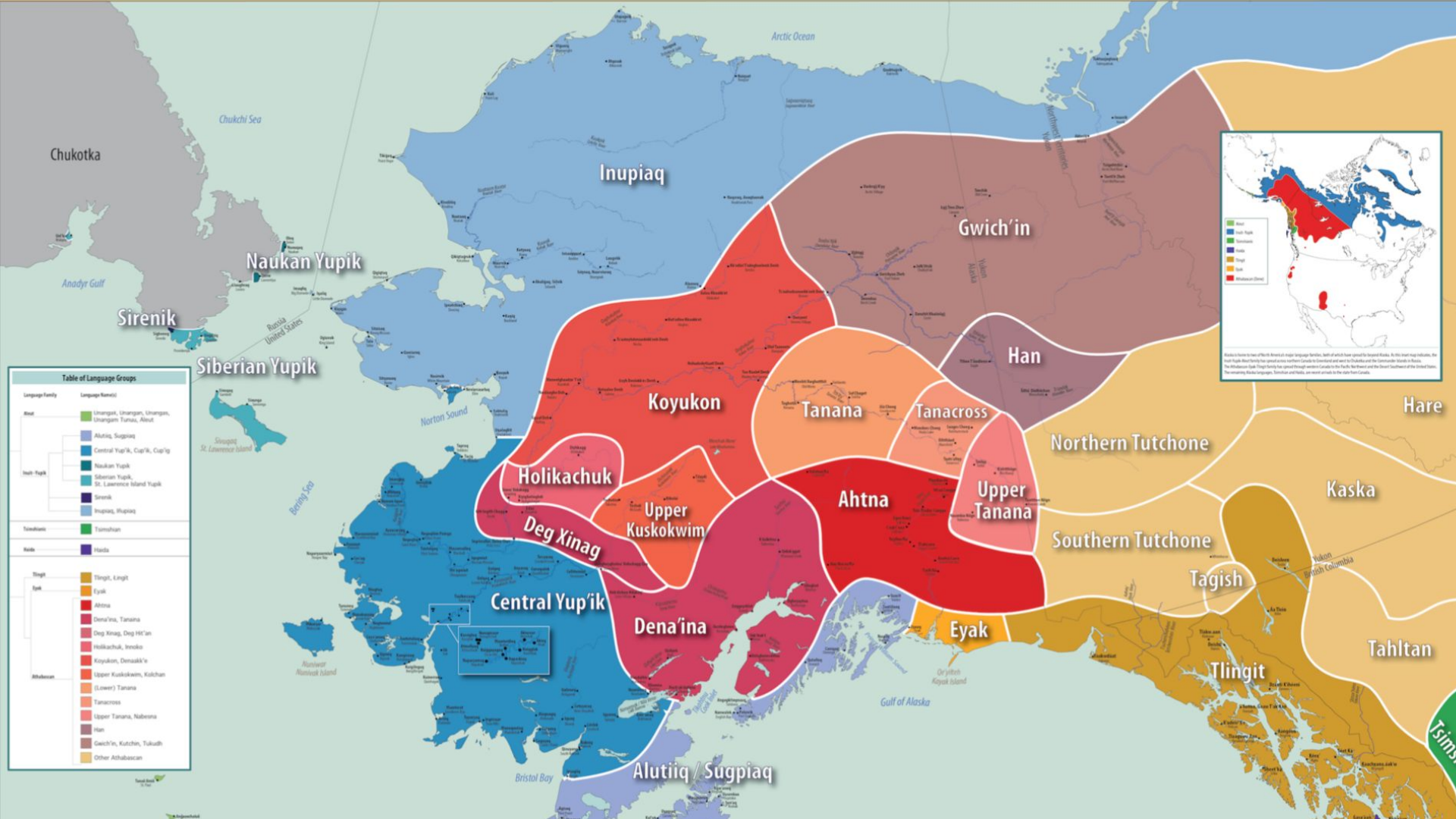


2019 | INTERNATIONAL YEAR OF
Indigenous Languages

- Increasing understanding, reconciliation and international cooperation.
- Creation of favorable conditions for knowledge-sharing & dissemination of good practices with regards to indigenous languages.
- Integration of indigenous languages into standard setting.
- **Empowerment through capacity building.**
- Growth and development through elaboration of new knowledge.







Alaska is home to over 100 Native languages. Major language families, both of which have spread to beyond Alaska, are the Eskimo-Aleut and the Athabaskan. The Eskimo-Aleut language family is spread throughout northern Canada to the Pacific Northwest and the coast of the United States. The Athabaskan language family is spread throughout western Canada to the Pacific Northwest and the coast of the United States. The remaining 80+ languages, including the Tlingit and Haida, are most closely related to the north-west coast.

Table of Language Groups

Language Family	Language Name(s)
Eskimo-Aleut	Unangan, Unangam, Unangam, Unangam Tunuu, Aleut
	Aleutic, Sugpiaq
	Central Yup'ik, Cap'ik, Cap'ig
	Naukan Yupik
	Siberian Yupik, St. Lawrence Island Yupik
Inuit-Yupik	Sirenik
	Inupiaq, Rupiaq
Tsimshian	Tsimshian
	Haida
Tlingit	Tlingit, Kligit
	Eyak
	Ahtna
	Chena'ina, Tanana
	Deg Xinag, Deg H'et'an
	Holikachuk, Inelko
	Koyukon, Demakik'ic
	Upper Kuskokwim, Kulkhan
	(Lower) Tanana
	Tanacross
Upper Tanana, Nabesna	
Athabaskan	Han
	Gwich'in, Kutchin, Tuluath
Other Athabaskan	

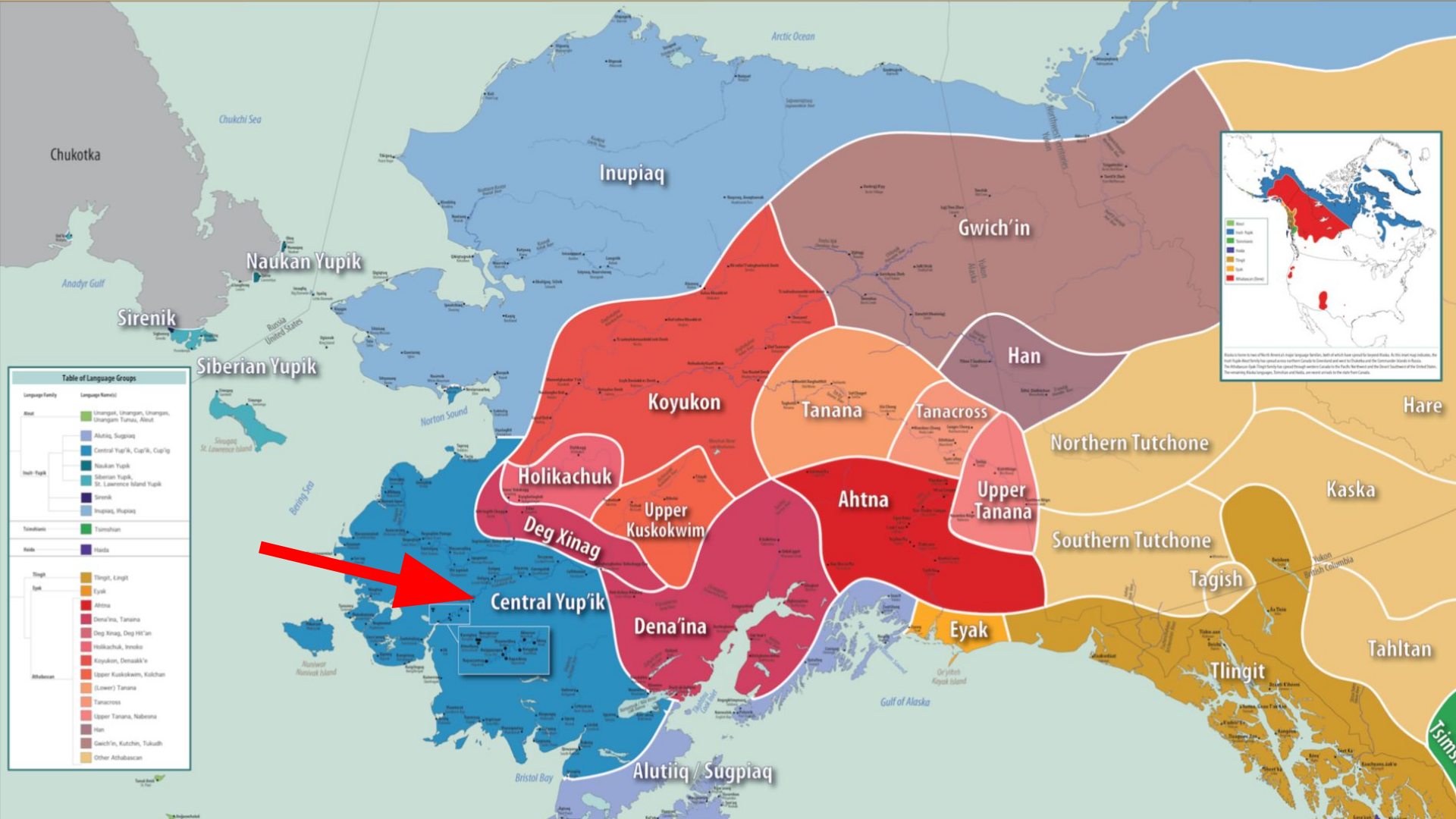
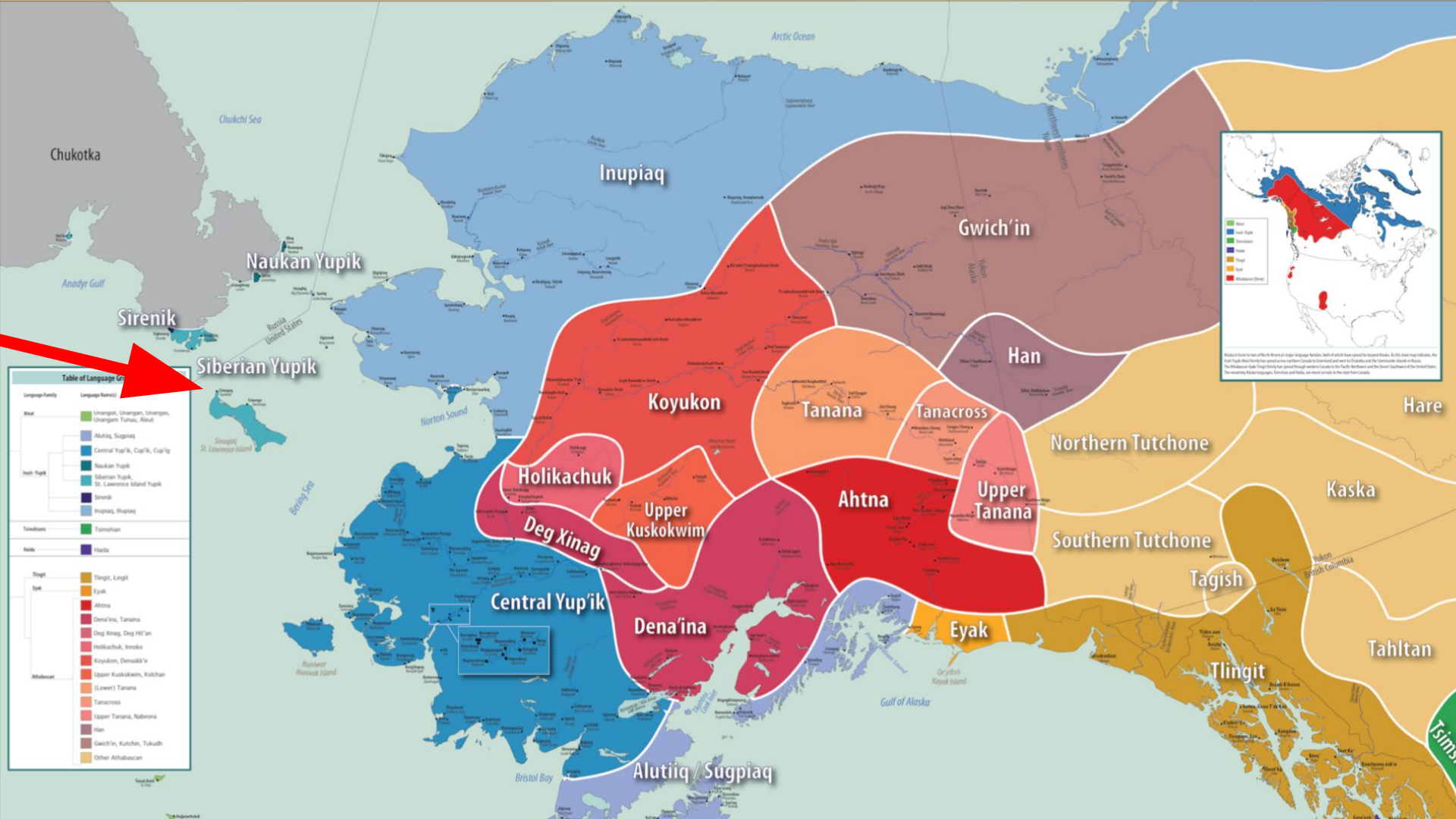


Table of Language Groups

Language Family	Language Name(s)
Eskimo	Unangan, Unangan, Unangan, Unangan Tunuu, Aleut
	Aleutic, Sugpiaq
	Central Yup'ik, Cap'ik, Cap'ig
	Naukan Yupik
	Siberian Yupik, St. Lawrence Island Yupik
Inuit-Yupik	Sirenik
	Inupiaq, Rupiaq
Tsimshian	Tsimshian
Haida	Haida
Athabaskan	Tlingit, Kligit
	Eyak
	Ahtna
	Chena'ina, Tanana
	Deg Xinag, Deg H'et'an
	Holikachuk, Ineeko
	Koyukon, Demakik'ic
	Upper Kuskokwim, Kulkhan
	(Lower) Tanana
	Tanacross
Upper Tanana, Nabesna	
Gwich'in	Gwich'in, Kutchin, Tukulth
	Other Athabaskan



Alaska is home to over 100 different languages. The map shows the distribution of major language families. Each of the colors represents a language family. The map is color-coded by language family. The map shows the distribution of major language families. The map is color-coded by language family. The map shows the distribution of major language families.



Alaska is home to over 100 different languages, many of which are endangered. The map shows the major language families and dialects in Alaska. The inset map shows the location of Alaska within North America. The Athabaskan language family is the largest and most diverse, followed by Eskimo-Aleut and Na-Dene. The map is color-coded by language family: Athabaskan (red), Eskimo-Aleut (blue), and Na-Dene (yellow). Other languages are shown in grey.

Table of Language Families

Language Family	Language Name(s)
Eskimo-Aleut	Unangan, Unangan, Unangan, Unangan Tunuu, Aleut
	Aleutic, Sugpiaq
	Central Yup'ik, Cap'ik, Cap'ig
	Naukan Yupik
	Siberian Yupik, St. Lawrence Island Yupik
Na-Dene	Sirenik
	Inupiaq, Rupiaq
Tutchone	Tutchone
	Tahltan
Athabaskan	Dena'ina
	Other Athabaskan





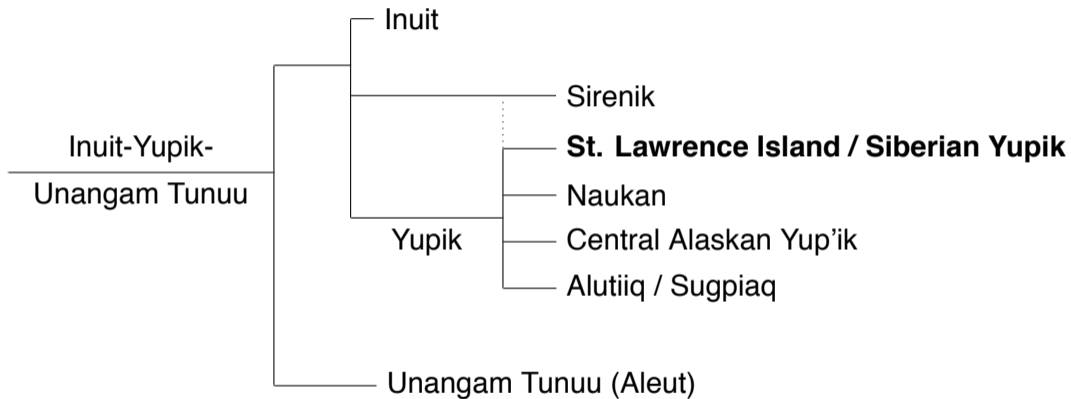




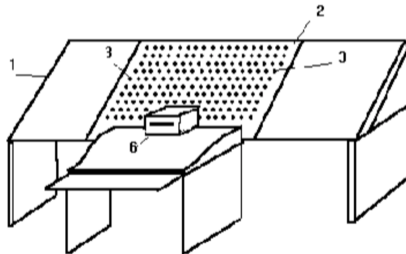
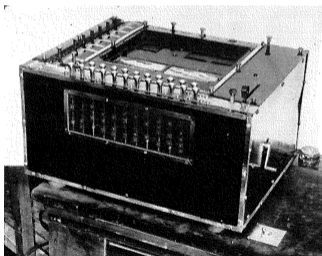
Inuit-Yupik-Unangam Tunuu language family

- Greenland (Inuit)
- Northern Canada (Inuit)
- Northern Alaska (Inuit)
- Western Alaska (Yup'ik)
- Southwestern Alaska (Sugpiak, Unangam Tunuu)
- **St. Lawrence Island (Yupik)**
- Big Diomedes (Inuit)
- Far eastern Russia (Yupik, Sirenik)

Inuit-Yupik-Unangam Tunuu language family

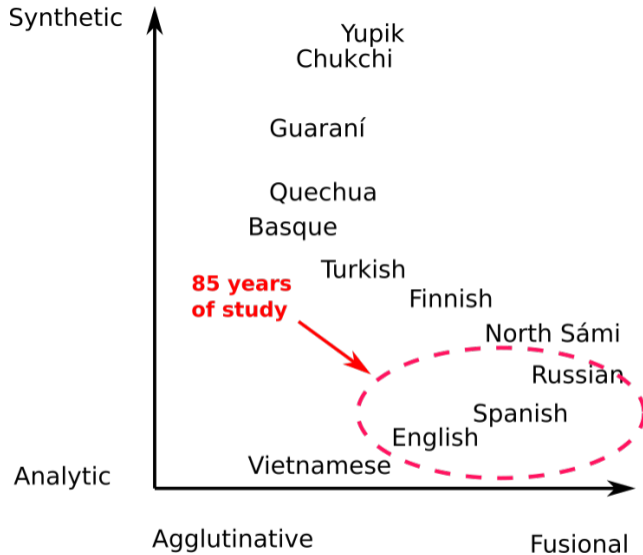


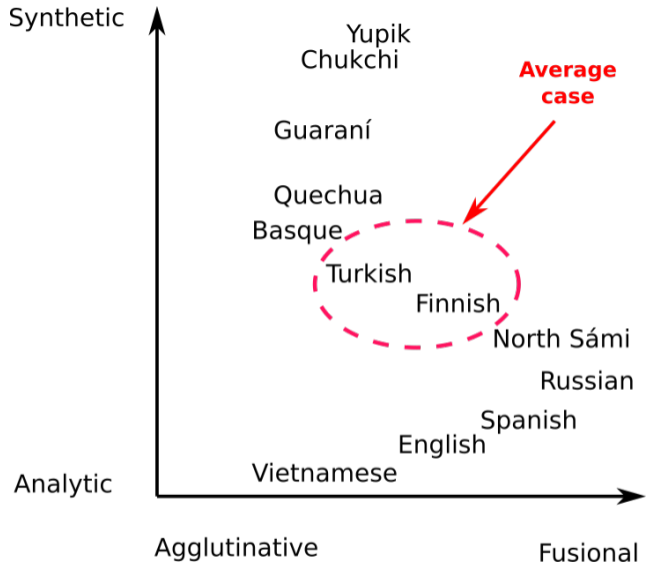
Since 1933, NLP technology has overwhelmingly focused on languages & methodologies in which the word is the primary meaning-bearing unit

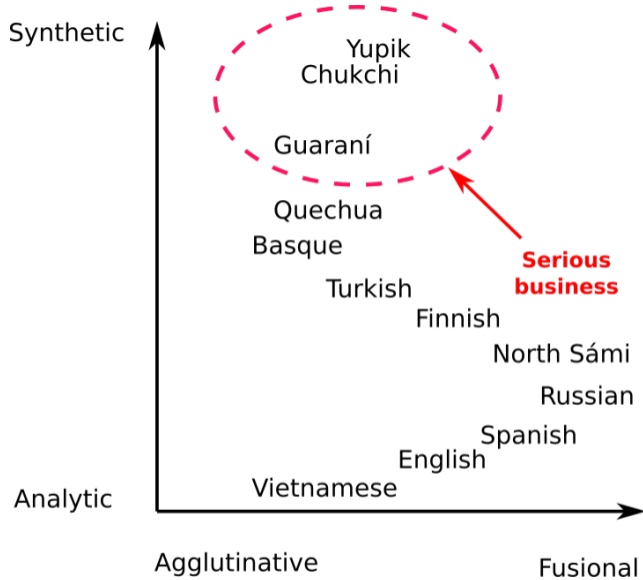


For *most* human languages, this assumption is **fundamentally broken**









$$p(\tau_t | \tau_1 \cdot \cdot \cdot \tau_{t-1})$$




$$p(\tau_t | \tau_1 \dots \tau_{t-1}) = \frac{\text{count}(\tau_1 \dots \tau_t)}{\text{count}(\tau_1 \dots \tau_{t-1})}$$





English



Yupik

* actual data disparity
is much much larger

dog

dogs

qikmiq

qikmik

qikmit

qikmiq

qikmik

qikmit

qikmigka

qikmigken

qikminka

qikmii

qikmikek

qikmiik

qikmiqa

qikmighpung

qikmighput

There are 1.2×10^{23}
stars in the observable universe.



There are 1.2×10^{23}
possible Yupik word forms.



Big data is NOT the solution.



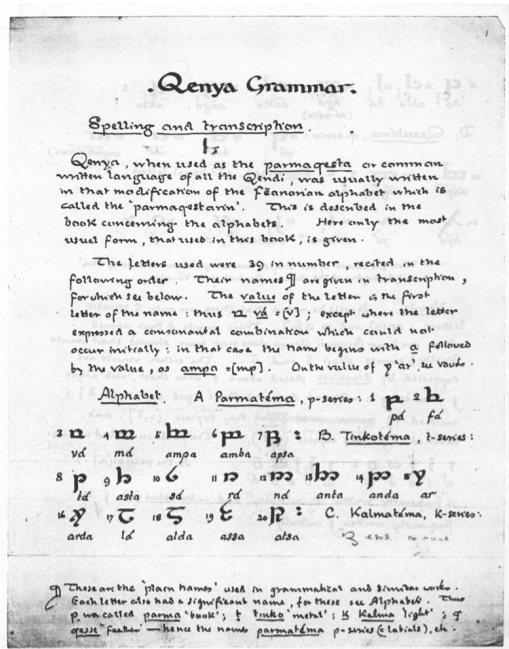
Modelling only at the word-level is like modelling only at a galaxy-level.





Հճառարարք՝ Գ. Գառարք՝
Գառարք՝ Գ. Գառարք՝





PARMA ELDALACBERON XXII

The Feanorian Alphabet · PART I

Quenya Verb STRUCTURE

by J. R. R. TOLKIEN

*Շահազարեայր՝է. Շահազար՝է
 Շահագիւղայր՝է. Կրնեմ շիգոյր՝է*

+



==

ash	nazg	durb-	-at-	-ul-	-ûk
one	ring	to.dominate	[Ptcp]	[3Pl]	[Compl]

...

agh	burzum	ishi	krimp-	-at-	-ul
and	darkness	inside	to.bind	[Ptcp]	[3Pl]

Course goals

- Learn about a new language from a reference grammar
 - Demonstrate your understanding through writing and teaching
- Select a topic from computational linguistics applicable to this language
 - Conduct a literature review, resulting in an annotated bibliography & report on state of the art
- Perform research on this topic
 - Identify state-of-the-art baseline, implement & extend it, run experiments, write a paper
- Conduct extended research in a group
 - Collaborate, experiment, and jointly author a paper
- Act as a peer reviewer for your classmates' work

- 1300 Yupiget on St. Lawrence Island
- 800 Yupiget on Russian mainland
- 300-400 Yupiget on Alaskan mainland

- 1930s-1950s Yupik materials developed in Russia
- 1970s-1990s Yupik materials developed in Alaska

- By mid-20th century, shift away from Yupik in Russia
- Current estimate of < 200 L1 Yupik speakers in Russia
- Youngest L1 Russian Yupiget estimated age > 70

Language shift - St. Lawrence Island

- In 1980, nearly all St. Lawrence Island Yupiget children spoke Yupik at home
- By mid-1990s through early 2000s, shift away from Yupik among SLI youth
- All SLI Yupiget born 1980 or earlier assumed to be L1 Yupik
- Current estimate of at least 540 L1 Yupik speakers on SLI
- Youngest L1 SLI Yupiget not known

Phonology & Orthography

Close Vowels	i i и	u u y	Latin IPA Cyrillic
Mid Vowel		e э ы	Latin IPA Cyrillic
Open Vowel	a a а		Latin IPA Cyrillic

Syllable structure

- Word-initial V(C)
- Otherwise CV(C)

- V may be short (e, a, i, u) or long (aa, ii, uu)
- Adjacent consonants only at syllable boundaries
- Adjacent consonant generally must agree in voicing

Phonology & Orthography

	Labial	Alveolar	Palatal	Retroflex	Velar	Velar (rounded)	Uvular	Uvular (rounded)	Glottal	
Unvoiced Stops	p p п	t t т			k k к	kw k ^w кӱ	q q к	qw q ^w кӱ		Latin IPA Cyrillic
Voiced Continuants	v v в	l l л	z z з	y j й	r ɻ р	g ɣ г	w ɣ ^w (r)ӱ	gh ɣ г	ghw ɣ ^w гӱ	Latin IPA Cyrillic
Unvoiced Continuants	f f ф	ll ɬ ль	s s с	rr ʂ ш	gg x х	wh x ^w xӱ	ghh χ х	ghhw χ ^w xӱ	h h г	Latin IPA Cyrillic
Voiced Nasals	m m м	n n н			ng ŋ ң	ngw ŋ ^w ңӱ				Latin IPA Cyrillic
Unvoiced Nasals	mm m̥ мь	nn n̥ нь			ngng ŋ̥ ңь	ngngw ŋ̥ ^w ңьӱ				Latin IPA Cyrillic

Legacy Digitization

Background

○○○○○

Intersecting machine learning & linguistic fieldwork

St. Lawrence Island Yupik

○○○○○○○○○○○○○○○○○○○○●○

Intersecting machine learning & linguistic fieldwork

○○○○○○○○○○

Moving forward

○○○○○

Lane Schwartz

- 3-volume Lore of St. Lawrence Island
- 3-volume Elementary Yupik readers
- 1-volume of Russian Yupik stories

Intersecting machine learning & linguistic fieldwork

- Yupik is polysynthetic, allowing for morphologically-complex words

(1) **mangteghaghllangllaghyugtukut**

mangteghagh-	-ghllag-	-ngllagh-	-yug-	-tu-	-kut
house-	-big-	-build-	-want.to-	-INTR.IND-	-1PL

'We want to build a big house'

- Yupik words typically adhere to the following template:

Root + 0-7 Derivational Morpheme(s) + Inflectional Morphemes + (Enclitic)

- Yupik is polysynthetic, allowing for morphologically-complex words

(1) **mangteghaghllangllaghyugtukut**

mangteghagh-	-ghllag-	-ngllagh-	-yug-	-tu-	-kut
house-	-big-	-build-	-want.to-	-INTR.IND-	-1PL

'We want to build a big house'

- Yupik words typically adhere to the following template:

Root + 0-7 Derivational Morpheme(s) + Inflectional Morphemes + (Enclitic)

- Yupik is polysynthetic, allowing for morphologically-complex words

(1) **mangteghaghllangllaghyugtukut**

mangteghagh-	-ghllag-	-ngllagh-	-yug-	-tu-	-kut
house-	-big-	-build-	-want.to-	-INTR.IND-	-1PL

'We want to build a big house'

- Yupik words typically adhere to the following template:

Root + 0-7 Derivational Morpheme(s) + Inflectional Morphemes + (Enclitic)

- Yupik is polysynthetic, allowing for morphologically-complex words

(1) **mangteghaghllangllaghyugtukut**

mangteghagh-	-ghllag-	-ngllagh-	-yug-	-tu-	-kut
house-	-big-	-build-	-want.to-	-INTR.IND-	-1PL

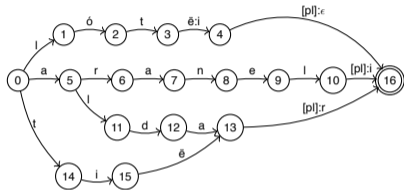
'We want to build a big house'

- Yupik words typically adhere to the following template:

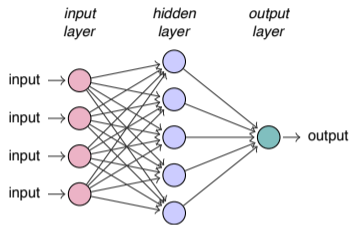
Root + 0-7 Derivational Morpheme(s) + **Inflectional Morphemes** + (Enclitic)

- Morphological analyzers may be implemented as a

Finite-State Transducer



Neural Network



- Neural systems require LOTS of data
 - But Yupik is a low-resource language
 - Very few surface form-lexical form pairs available

- **OBJECTIVE:** Analyze inflected Yupik nouns with no derivational morphology
- **TRAINING DATA:** Every nominal surface form and its respective lexical form
 - 3873 Yupik noun roots
 - 273 inflectional suffixes
 - $3873 \times 273 = 1,057,329$ total nouns
 - 658,410 after removing duplicate surface forms (case syncretism)

Surface Form	Lexical Form
mangteghaq	mangteghagh[N][ABS][SG]
mangteghaat	mangteghagh[N][ABS][PL]
mangteghaak	mangteghagh[N][ABS][DU]
mangteghaa	mangteghagh[N][ABS][SG][3SGPOSS]
⋮	

- **EVALUATION OBJECTIVES**

- Evaluate on a neutral dataset
- Contrast performance with the FST analyzer

- **NEUTRAL DATASET:** *Mrs. Della Waghiyi's St. Lawrence Island Yupik Texts With Grammatical Analysis* (Waghiyi & Nagai, 2001)

- Identified **344 inflected nouns with no derivational morphology**

- Supplemented the FST analyzer with a guesser module

- Results:

	Coverage (%)	Accuracy (%)
FST (No Guesser)	85.96	79.82
FST (w/Guesser)	100	84.50
Neural	100	91.81

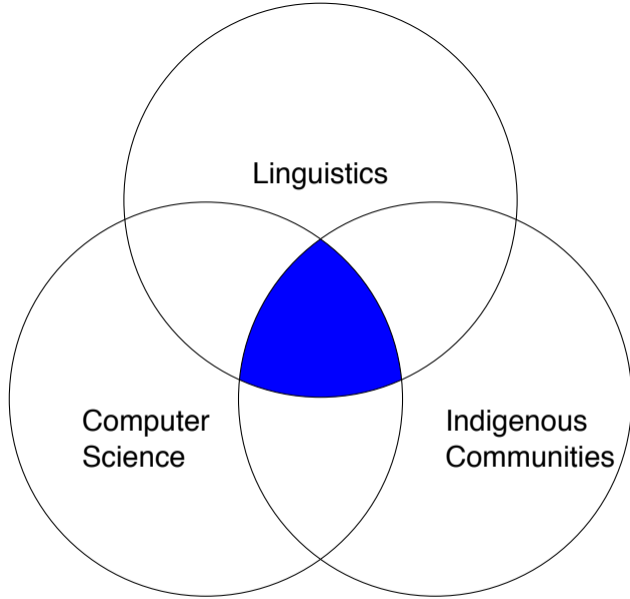
- An **out-of-vocabulary** (OOV) root is an unattested root that appears in the Waghiyi & Nagai (2001) evaluation dataset but does not appear in our data

OOV Root	FST	NN
aghnasinghagh	–	–
aghveghniigh	–	✓
akughvigagh	✓	✓
qikmiraagh	–	–
sakara	✓	–
sanaghte	–	–
tangiqagh	–	✓

- A root with a **spelling variant** is one that differs in the Waghiyi & Nagai (2001) evaluation set from its form in our data

Root Variant	FST	NN
melqighagh	✓	✓
piitesiiighagh	–	✓
uqfiilleghagh	–	✓
*ukusumun	–	✓

Building a virtuous cycle



- Digitization of legacy materials
- Pedagogical materials & tools
- Orthographic experimentation
- Identify under-described phenomena
- Real-time morphological analysis

- Digitization of legacy materials
- **Pedagogical materials & tools**
- Orthographic experimentation
- Identify under-described phenomena
- Real-time morphological analysis

- Digitization of legacy materials
- Pedagogical materials & tools
- **Orthographic experimentation**
- Identify under-described phenomena
- Real-time morphological analysis

- Digitization of legacy materials
- Pedagogical materials & tools
- Orthographic experimentation
- **Identify under-described phenomena**
- Real-time morphological analysis

- Digitization of legacy materials
- Pedagogical materials & tools
- Orthographic experimentation
- Identify under-described phenomena
- **Real-time morphological analysis**

Moving forward

$$p(\mathbf{e}) = p(e_t | e_1 \dots e_{t-1})$$

$$p(\mathbf{e}) = p(e_t | e_1 \dots e_{t-1})$$
$$\approx p(e_t | e_{t-1})$$

1. Zero-order approximation

XFOML RXKHRJFFJUJ ZLPWCFWKCYJ
FFJEYVKCQSGXYD QPAAMKBZAACIBZLHJQD

2. First-order approximation

OCRO HLI RGWR NMIELWIS EU LL NBNESEBYA TH EEI
ALHENHTTPA OOBTTVA NAH BRL

3. Second-order approximation

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY
ACHIN D ILONASIVE TUCOOWE AT TEASONARE FUSO
TIZIN ANDY TOBE SEACE CTISBE

5. First-Order Word Approximation

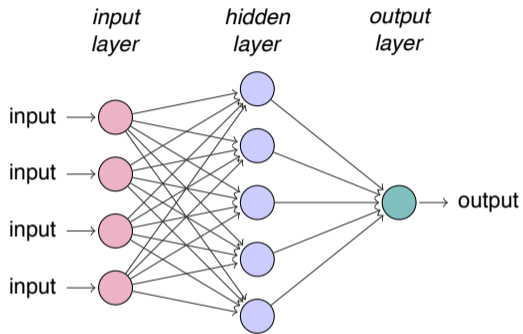
REPRESENTING AND SPEEDILY IS AN GOOD APT OR
COME CAN DIFFERENT NATURAL HERE HE THE A IN
CAME THE TO OF TO EXPERT GRAY COME TO FUR-
NISHES THE LINE MESSAGE HAD BE THESE.

6. Second-Order Word Approximation

THE HEAD AND IN FRONTAL ATTACK ON AN ENGLISH
WRITER THAT THE CHARACTER OF THIS POINT IS
THEREFORE ANOTHER METHOD FOR THE LETTERS
THAT THE TIME OF WHO EVER TOLD THE PROBLEM
FOR AN UNEXPECTED

$$p(\mathbf{e}) = p(e_t | e_1 \dots e_{t-1})$$
$$\approx p(e_t | e_{t-1})$$

$$p(\mathbf{e}) = p(e_t | e_1 \dots e_{t-1})$$



- Legacy text digitization
- Web portal / interactive e-books
- App-based dictionary
- Language learning lessons
- foma-based spell-checker
- Forced aligner / speech recognizer
- Machine translation

Feature-rich Open-vocabulary Interpretable Language Modelling

Interpretable Tensor Morpheme Representation

yaławtəma
“with the head”

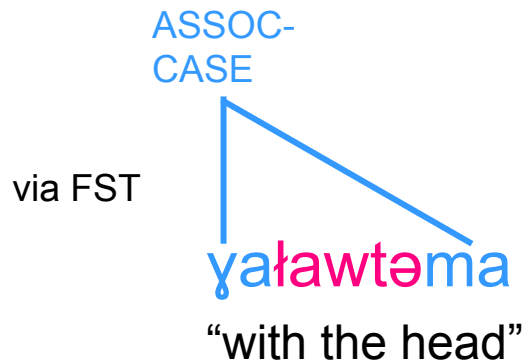
-*lewət*- “head” (Chukchi)

Interpretable Tensor Morpheme Representation

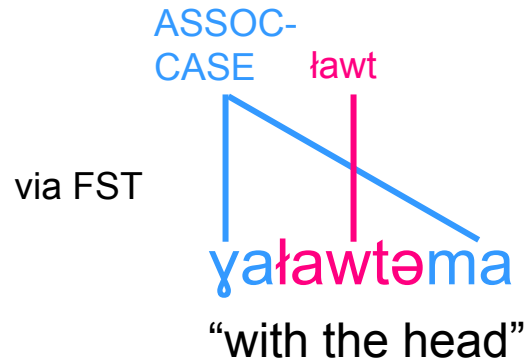
yaławtəma

“with the head”

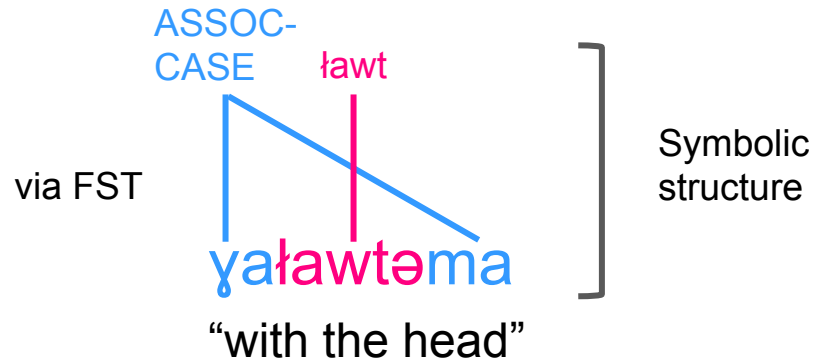
Interpretable Tensor Morpheme Representation



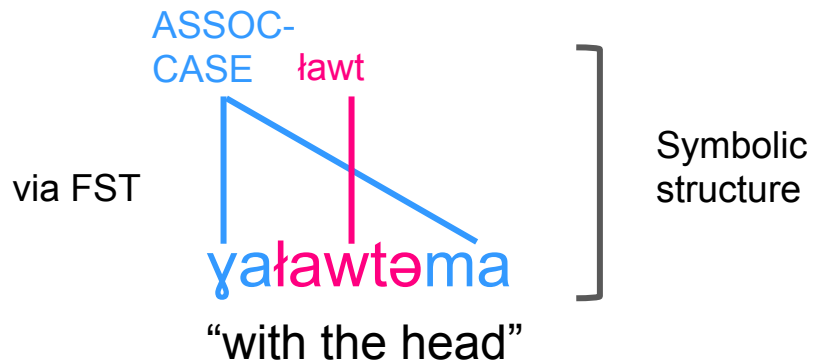
Interpretable Tensor Morpheme Representation



Interpretable Tensor Morpheme Representation

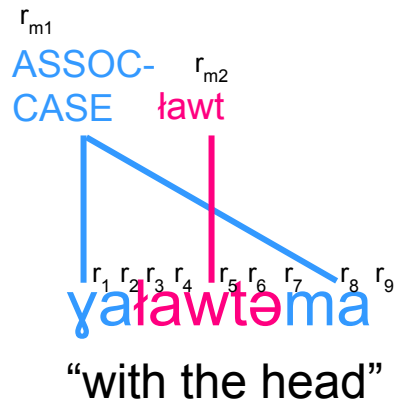


Interpretable Tensor Morpheme Representation



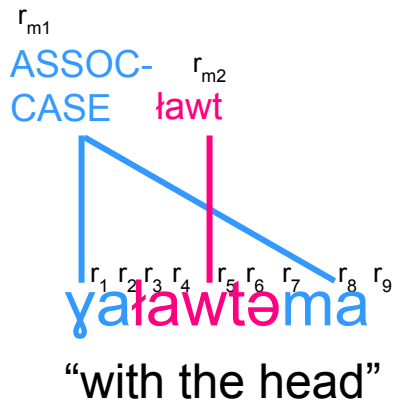
Symbolic structure → Tensor representation → Vector representation

Interpretable Tensor Morpheme Representation



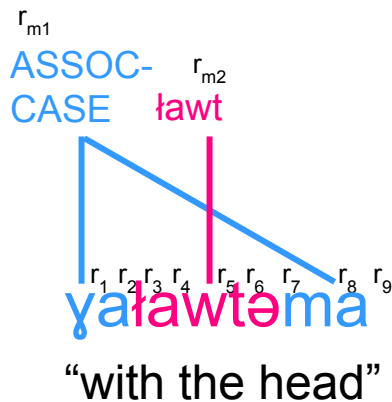
Decompose into *fillers* and *roles*.
(Smolenksy 1990)

Interpretable Tensor Morpheme Representation



Embed the fillers and roles into *vectors*
(Smolenksy 1990)

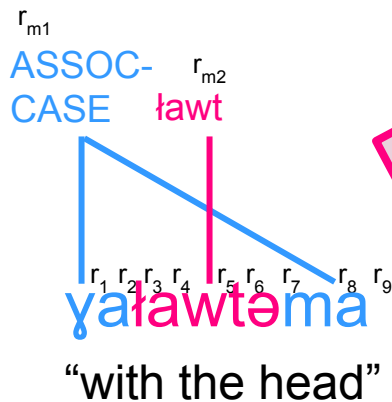
Interpretable Tensor Morpheme Representation



Embed the fillers and roles into *vectors*
(Smolenksy 1990)

$$Repr(\text{ya-ma}) = (\hat{y} \otimes \hat{r}_1 + \hat{a} \otimes \hat{r}_2 + \hat{m} \otimes \hat{r}_8 + \hat{a} \otimes \hat{r}_9) \otimes \hat{r}_{m1}$$

Tensor Morpheme Representation

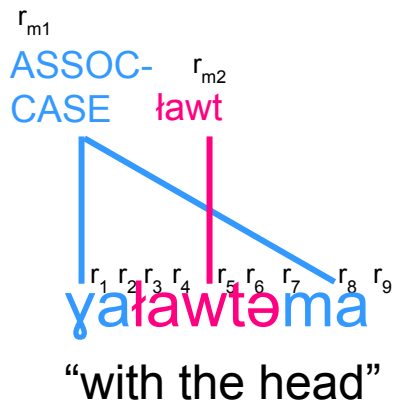


This means we have embeddings that are interpretable! $\hat{\text{y}}$ $\hat{\text{a}}$ $\hat{\text{m}}$

Embed the fillers and roles into *vectors*
(Smolenksy 1990)

$$\text{Repr}(\text{ya-ma}) = (\hat{\text{y}} \otimes \hat{r}_1 + \hat{\text{a}} \otimes \hat{r}_2 + \hat{\text{m}} \otimes \hat{r}_8 + \hat{\text{a}} \otimes \hat{r}_9) \otimes \hat{r}_{m_1}$$

Interpretable Tensor Morpheme Representation



1. Deterministically create these with FST for known sequences
2. Learn them with neural model (e.g. RNN seq2seq) to generalize

Embed the fillers and roles into *vectors*
(Smolensky 1990)

Deterministically construct morpheme tensors

- a. Run morphological analyzer on training data to identify morphemes

qikmighhaak \xrightarrow{a} qikmigh - ghhagh - [Abs.Du]
“Two small dogs” “dog - small.N - [Abs.Du]”

Deterministically construct morpheme tensors

- Run morphological analyzer on training data to identify morphemes
- Use **Tensor Product Representation** to deterministically calculate morpheme tensors



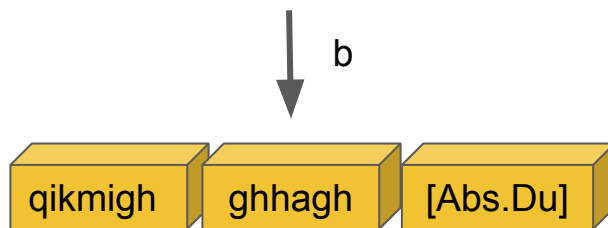
(St. Lawrence Island Yupik)

Deterministically construct morpheme tensors

- Run morphological analyzer on training data to identify morphemes
- Use **Tensor Product Representation** to deterministically calculate **morpheme tensors**
- Save these morpheme tensors for later use as gold standard labels

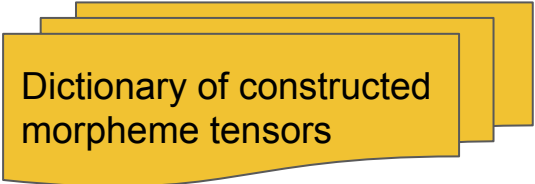
qikmighhaak \xrightarrow{a} qikmigh - ghhagh - [Abs.Du]
“Two small dogs” “dog - small.N - [Abs.Du]”

Dictionary of constructed morpheme tensors



(St. Lawrence Island Yupik)

Autoencoder



Dictionary of constructed
morpheme tensors

High dimensionality:
 $10^3 - 10^9$ floats per vector



qikmigh



ghhagh



[Abs.Du]

(St. Lawrence Island Yupik)

Dictionary of constructed
morpheme tensors

High dimensionality:
 $10^3 - 10^9$ floats per vector

qikmigh

ghhagh

[Abs.Du]

qikmigh

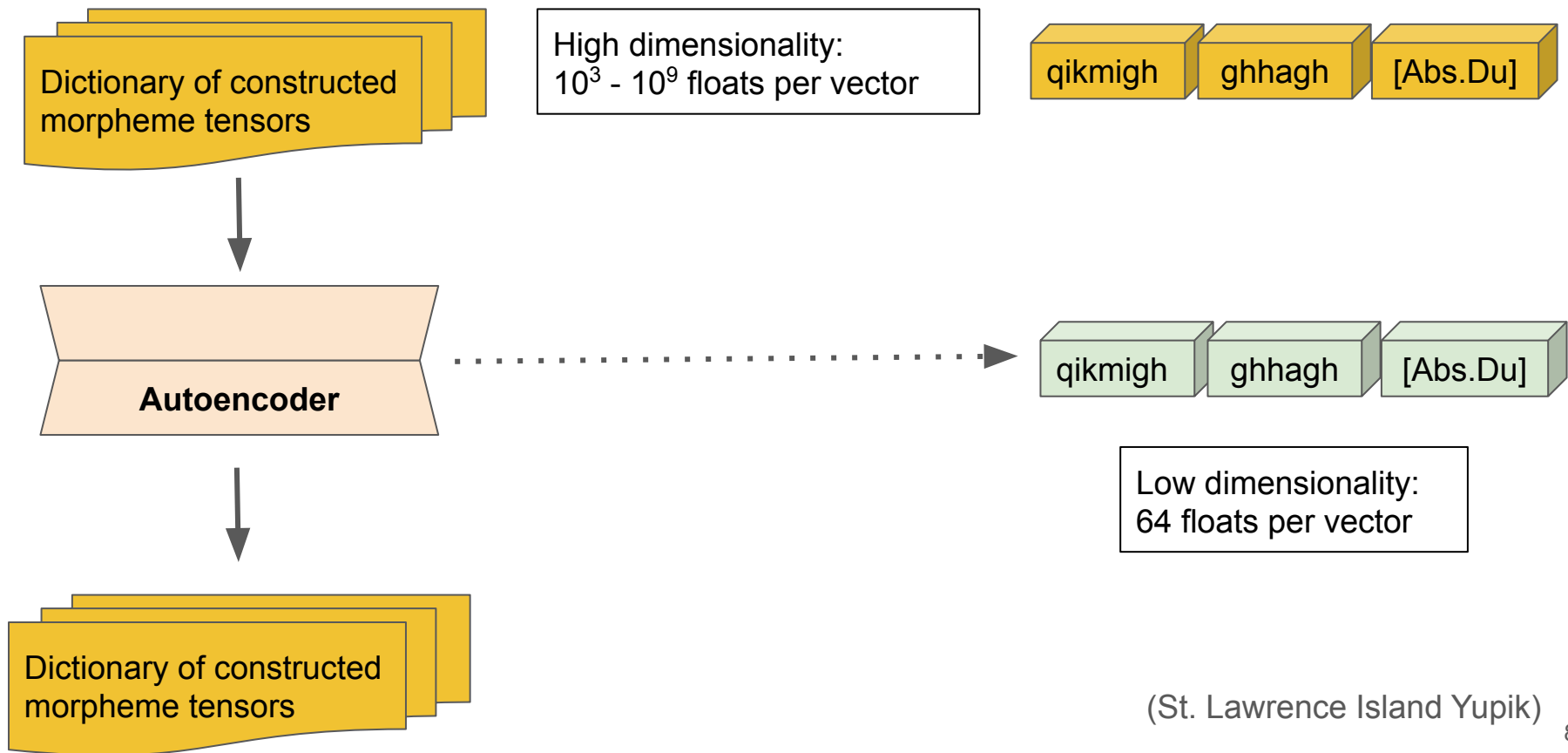
ghhagh

[Abs.Du]

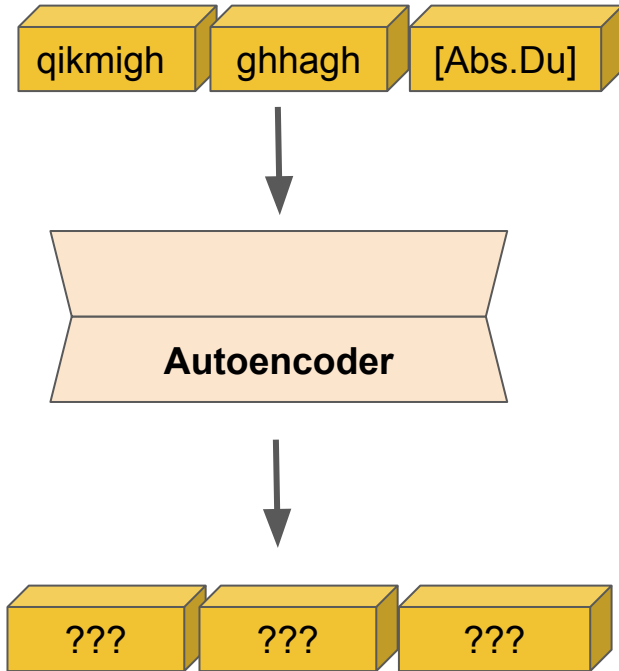
Low dimensionality:
64 floats per vector

(St. Lawrence Island Yupik)

Use autoencoder to learn morpheme vectors



Problem: Morpheme tensors are sparse



As a result, learning signal is very weak.

Solution: Unbinding Loss

