

Neural Polysynthetic Language Modelling

JSALT 2019 Opening Presentation











Indigenous Languages

matter for
development,
peace building
and reconciliation

Languages play a crucial role in the daily lives of people, not only as a tool for communication, education, social integration and development, but also as a repository for each person's unique identity, cultural history, traditions and memory. But despite their immense value, languages around the world continue to disappear at an alarming rate. With this in mind, the United Nations declared 2019 The Year of Indigenous Languages (IY2019) in order to raise awareness of them, not only to benefit the people who speak these languages, but also for others to appreciate the important contribution they make to our world's rich cultural diversity.





2019 | INTERNATIONAL YEAR OF
Indigenous Languages

7
thousand

—
Languages
spoken
worldwide

370
million

—
Indigenous
people
in the world

90
countries

—
With
indigenous
communities

5
thousand

—
Different
indigenous
cultures

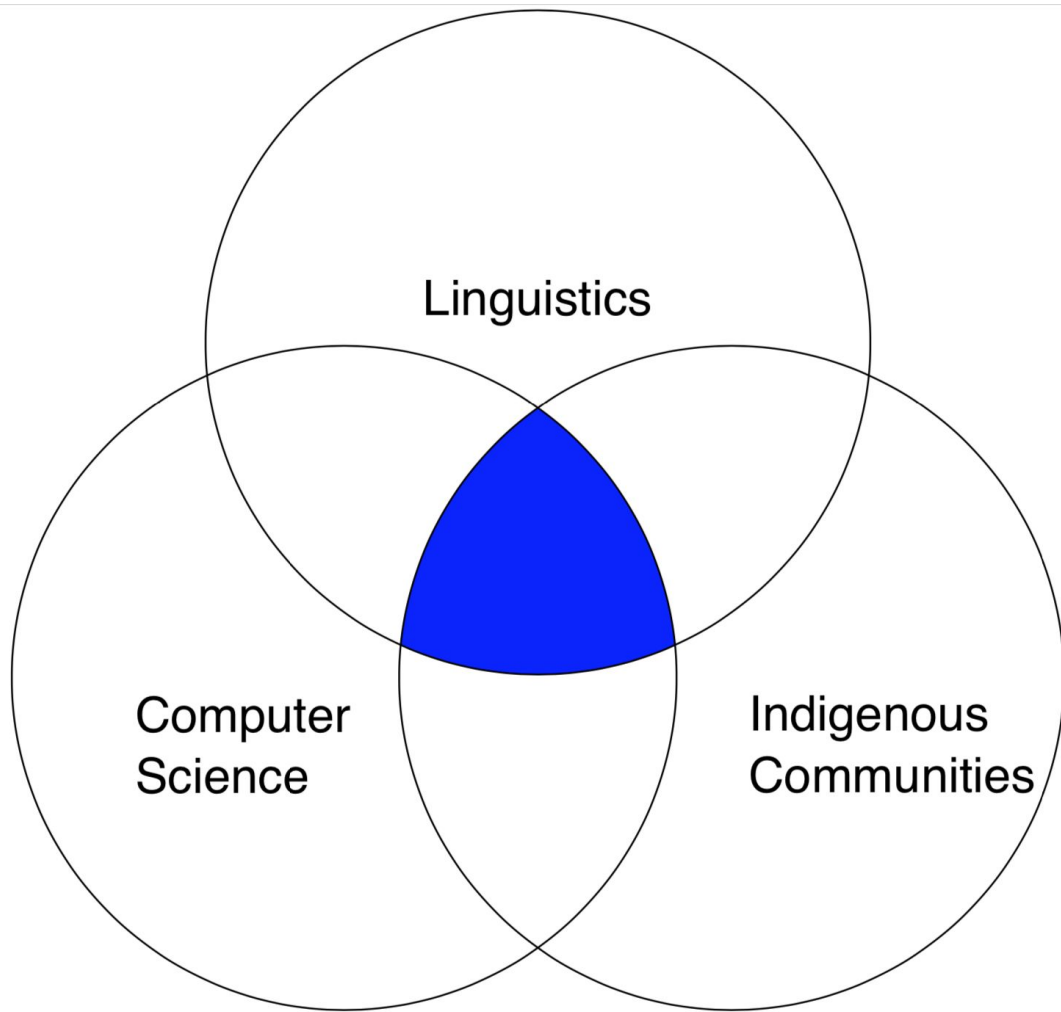
2680
languages

—
In danger

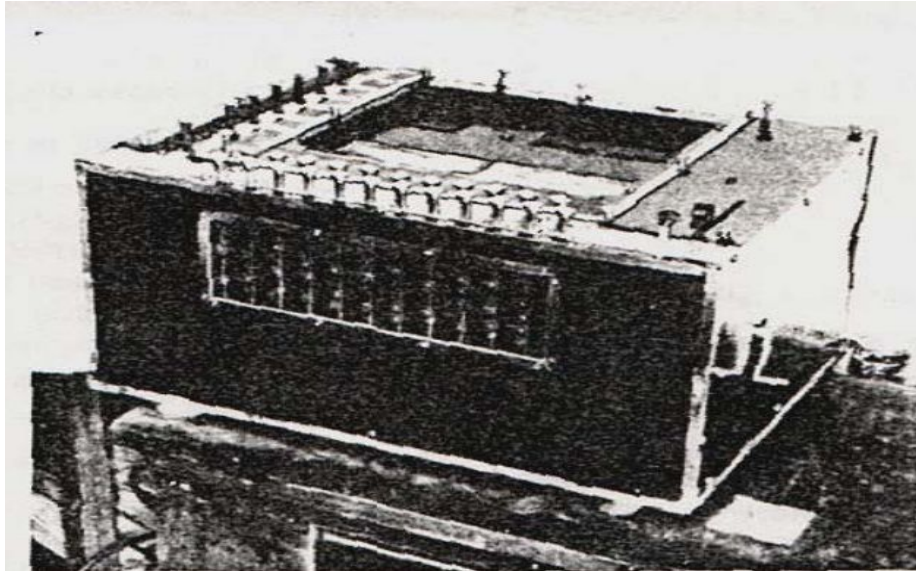


2019 | INTERNATIONAL YEAR OF
Indigenous Languages

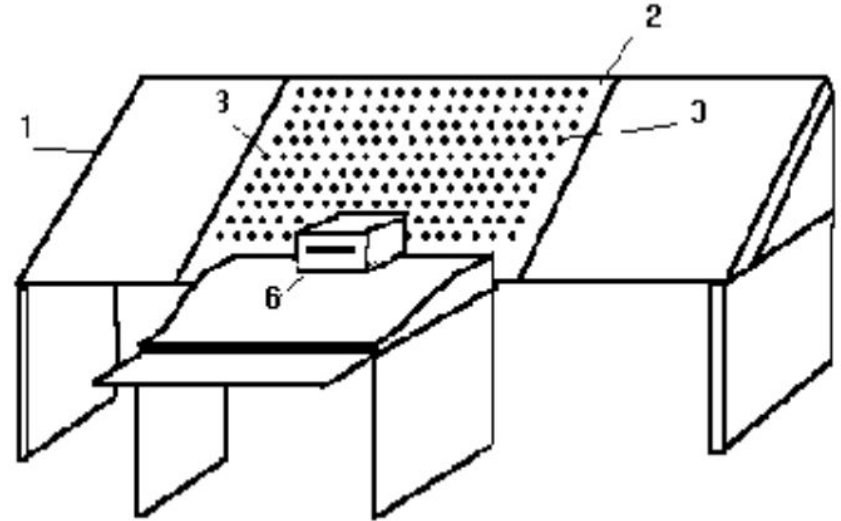
- Increasing understanding, reconciliation and international cooperation.
- Creation of favorable conditions for knowledge-sharing & dissemination of good practices with regards to indigenous languages.
- Integration of indigenous languages into standard setting.
- **Empowerment through capacity building.**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



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



Arstrouni (1933, Paris)



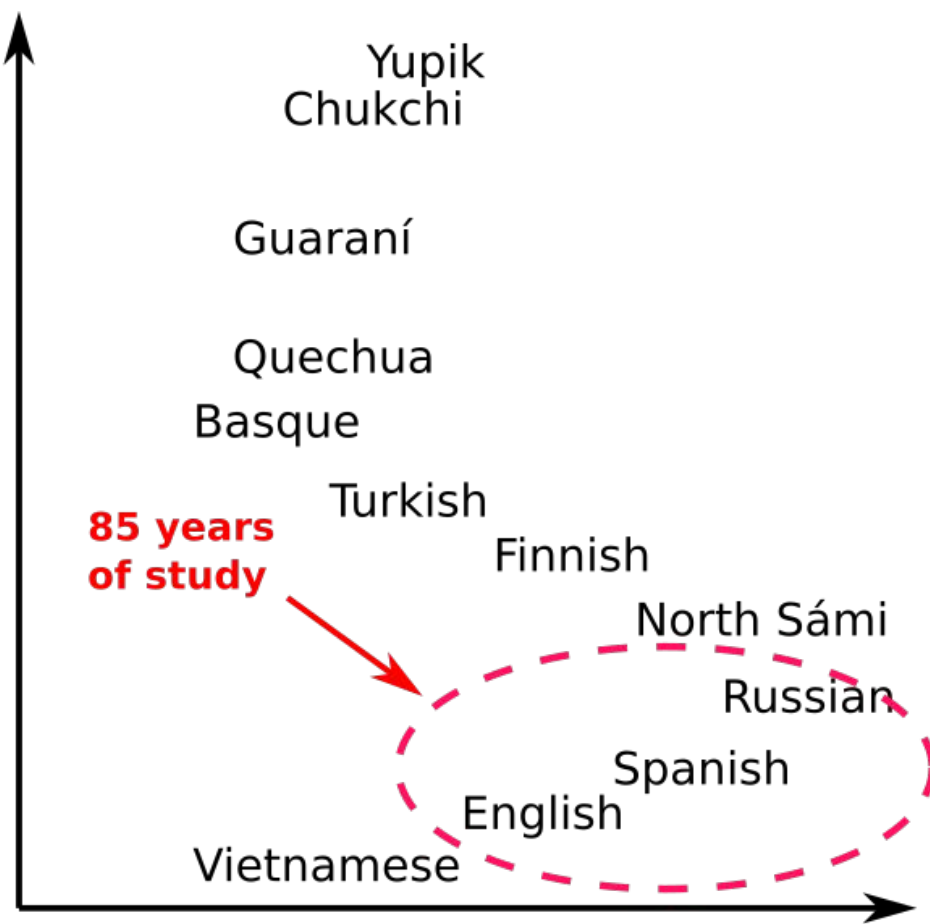
Trojanskij (1933, St. Petersburg)

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



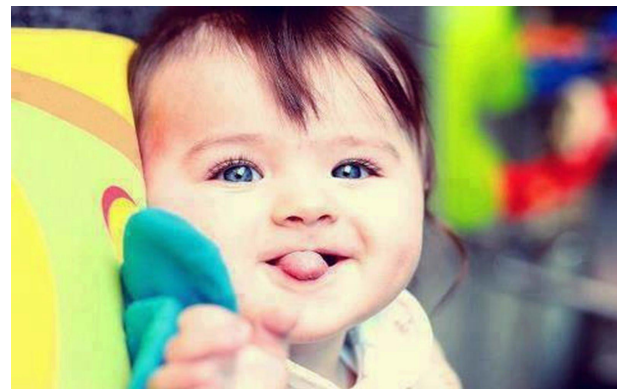
Synthetic

Analytic

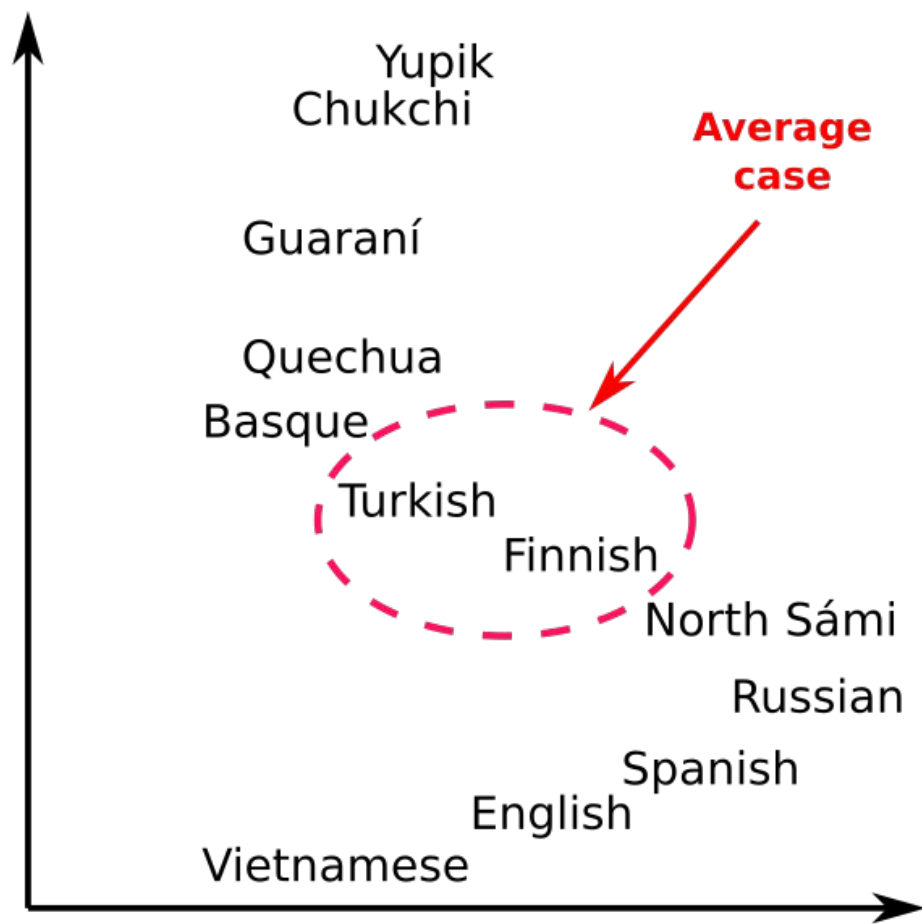


Agglutinative

Fusional



Synthetic

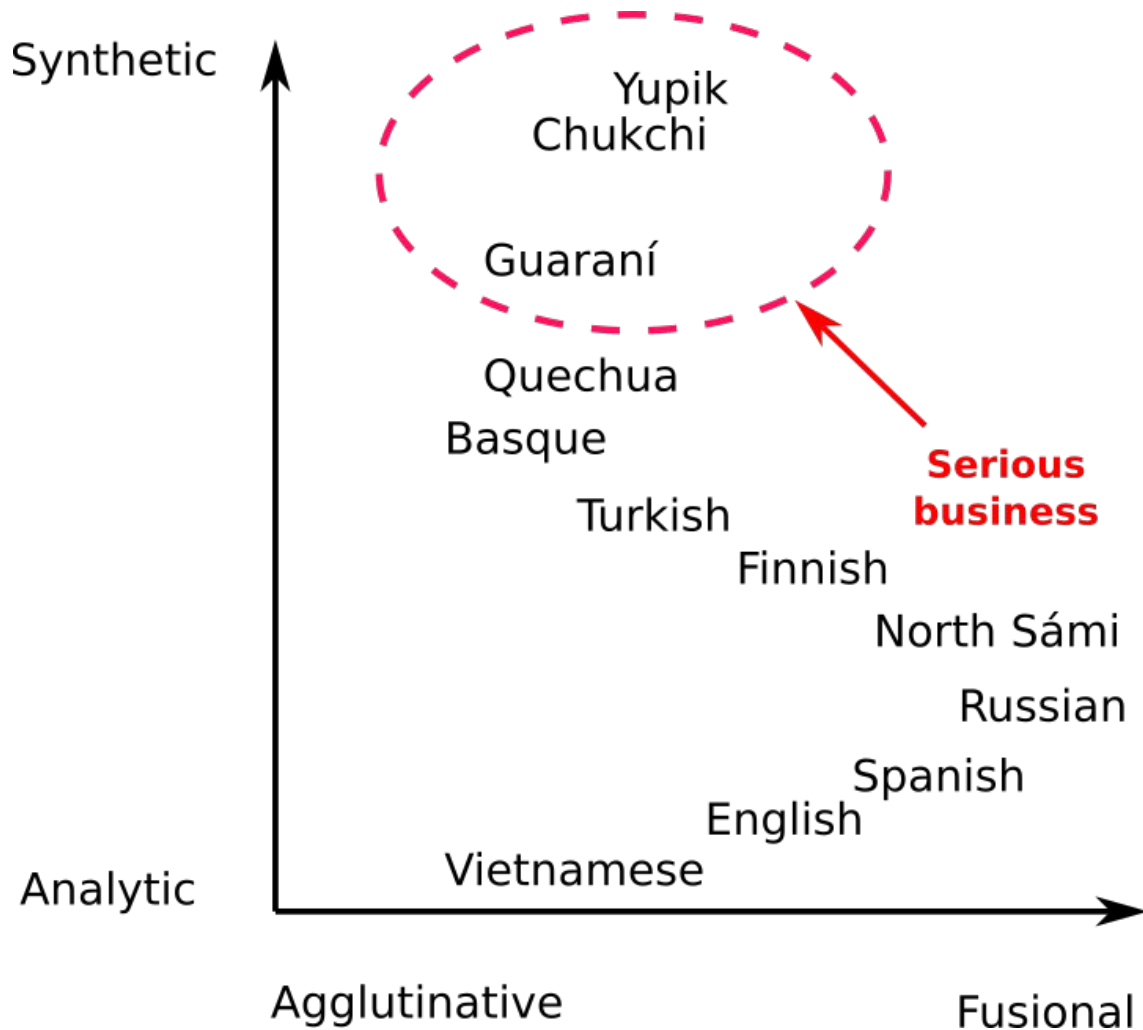


Analytic

Agglutinative

Fusional





Analytic Languages:

Little inflectional morphology

- **English** inflection
 - dog, dogs
 - walk, walks, walked, walking
- **English** derivation
 - disestablishmentarianism

Isolating Languages:

Little morphology other than compounding

- **Chinese** inflection
 - few affixes (prefixes and suffixes):
 - 们: 我们, 你们, 他们, 同志们
mén: *wǒmén*, *nǐmén*, *tāmén*, *tóngzhìmén*
plural: we, you (pl.), they, comrades, LGBT people
 - “suffixes” that mark aspect: 着 *-zhě* ‘continuous aspect’
- **Chinese** derivation
 - 艺术家 *yìshùjiā* ‘artist’
- **Chinese** is a champion in the realm of compounding—up to 80% of **Chinese** words are actually compounds.

毒	+	贩	→	毒贩
<i>dú</i>		<i>fàn</i>		<i>dúfàn</i>
‘poison, drug’		‘vendor’		‘drug trafficker’

Synthetic Languages:

Lots of morphology relative to analytic languages

- **Czech** inflection
 - fusional suffixes:

Class	Singular	Plural
1. Nominative	město	města
2. Genitive	města	měst
3. Dative	městu	městům
4. Accusative	město	města
5. Vocative	město	města
6. Locative	městě	městech
7. Instrumental	městem	městy

Agglutinative Languages:

Verbs in **Swahili** have an average of 4-5 morphemes, <http://wals.info/valuesets/22A-swa>

- Words are written without hyphens or spaces between morphemes.
- **Orange** prefixes mark noun class (like gender, except **Swahili** has nine instead of two or three).
 - Verbs agree with nouns in noun class.
 - Adjectives also agree with nouns.
 - Very helpful in parsing.
- **Black** prefixes indicate tense.

Swahili	English
<i>m-tu a-li-lal-a</i>	'The person slept'
<i>m-tu a-ta-lal-a</i>	'The person will sleep'
<i>wa-tu wa-li-lal-a</i>	'The people slept'
<i>wa-tu wa-ta-lal-a</i>	'The people will sleep'

Agglutinative Languages:

Turkish

But most words have around three morphemes

uygarlaştıramadıklarımızdanmışsınızcasına

“(behaving) as if you are among those whom we were not able to civilize”

uygar “civilized”

+laş “become”

+tır “cause to”

+ama “not able”

+dık past participle

+lar plural

+ımız first person plural possessive (“our”)

+dan ablative case (“from/among”)

+mış past

+sınız second person plural (“y’all”)

+casına finite verb → adverb (“as if”)

Derivational and inflectional morphology.

From the Jurafsky and Martin textbook.

Polysynthetic Languages

- All definitions are problematic
 - “A word means a whole sentence”
- There is no one feature (e.g., noun incorporation) that is in every polysynthetic language
- The definition isn't important
- What is important
 - There is a much larger number of distinct possible words
 - Beyond Turkish and Finnish
 - All languages have some of the properties of polysynthesis
 - Compounding, mild recursion in morphology (*operationalization*), causative, applicative, passive, negative

Inuit-Yupik Language Family

Yupik

St. Lawrence Island

Central Alaskan

Inuit

Inupiaq

Innuinaqtun

Inuktitut

Greenlandic (Kalaallisut)



Inuktitut Morphology

And even though it's not snowing a great deal, I'm not going out.

Qanniqlaunngikkaluaqtuqlu aninngittunga.

qanniq-+-lak + uq-+-nngit-+-galuaq-+-tuq + lu
to snow a little frequently not although 3rd pers. sg. and

ani-+-nngit-+-junga
to go out not 1st pers. sg.

Inuktitut Morphology

qanniq + lak+uq+nngit+galuaq + tuq + lu

ROOT	LEXICAL AFFIXES*	GRAMMATICAL ENDING	(CLITIC)
------	------------------	-----------------------	----------

	Zero or many (*)		
	Derivational – e.g. change v to n, n to v	Subject/object of verbs	
Verbal or Nominal	Semantic – adverbial, negative, tense	(fused with mood/mode)	small set, zero or 1
	Light verbs	Case	
	Certain adjectives	Possessors of nouns	
	incorporated into nouns	(fused with case)	

Inuktitut Morphology

Some examples of flipping back and forth between noun and verb

umiaq + juaq + liuq + vik + mi

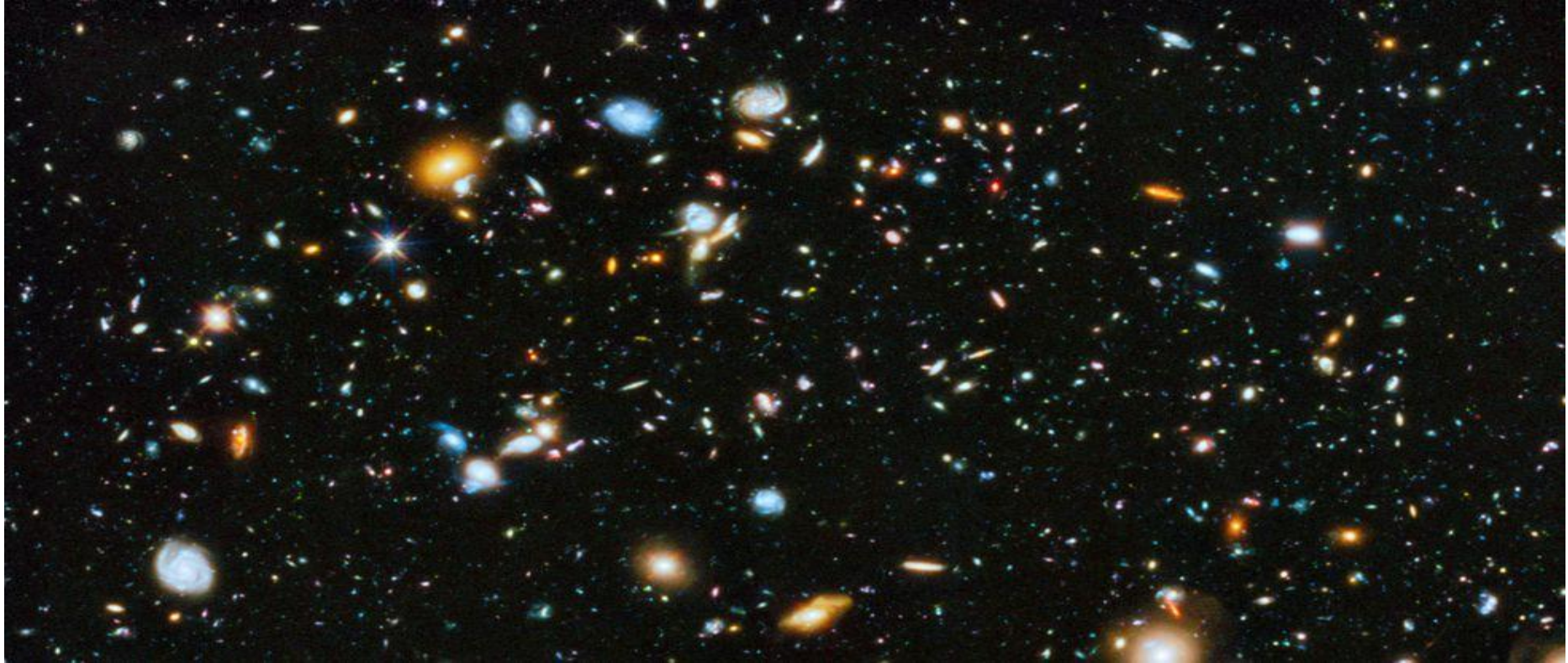
boat (n) + big (adj incorp) + make (n-v, light verb) + place-where (v-n, derivational) + locative
“in the shipyard”

ilinniaq + vik + siuq + junga

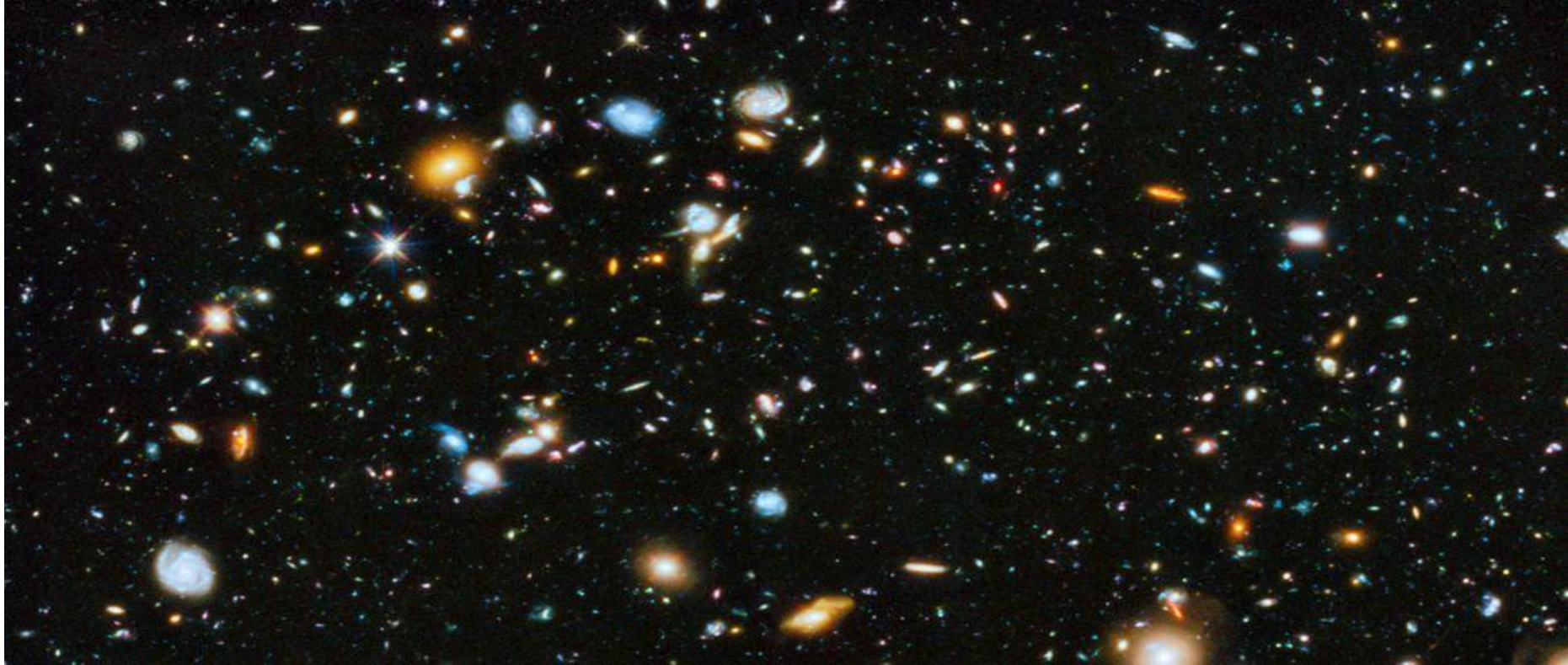
learn (v) + place-where (v-n, derivational) + look for (n-v, light verb) + 1-s
“I’m looking for a school”

<http://www.inuktitutcomputing.ca/Technocrats/ILFT.php#morphology>

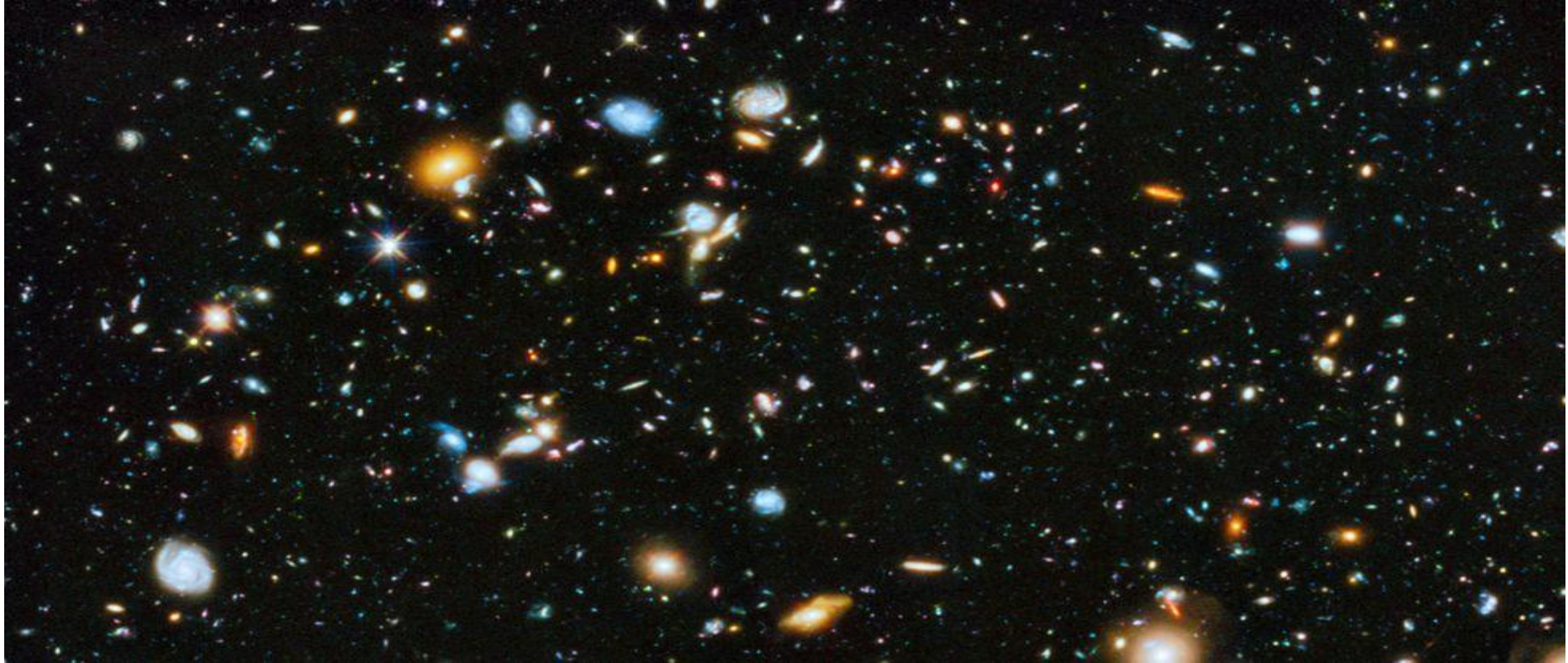
There are 1.2×10^{23} stars in the observable universe



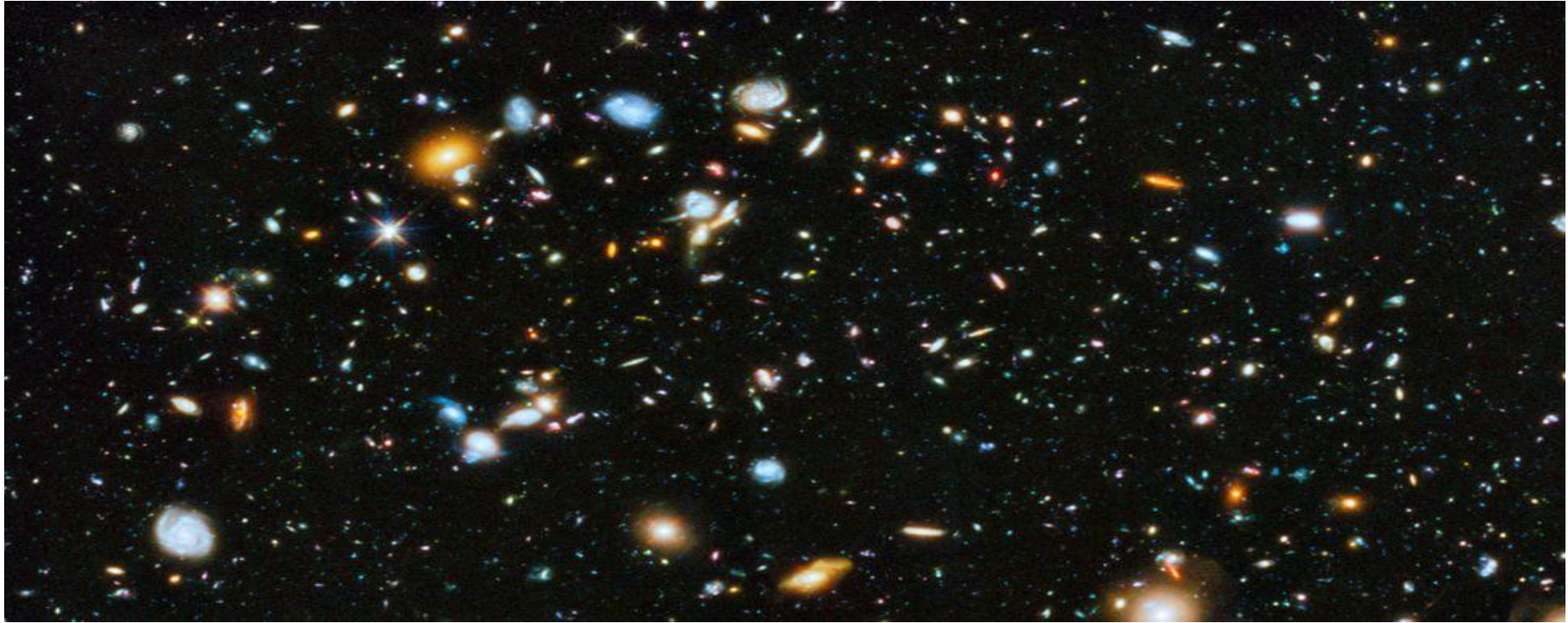
There are 1.2×10^{23} possible Yupik word forms



Big data is NOT the solution



We need a language model that sees the stars, not just the galaxies



Create a language model that handles the hardest case



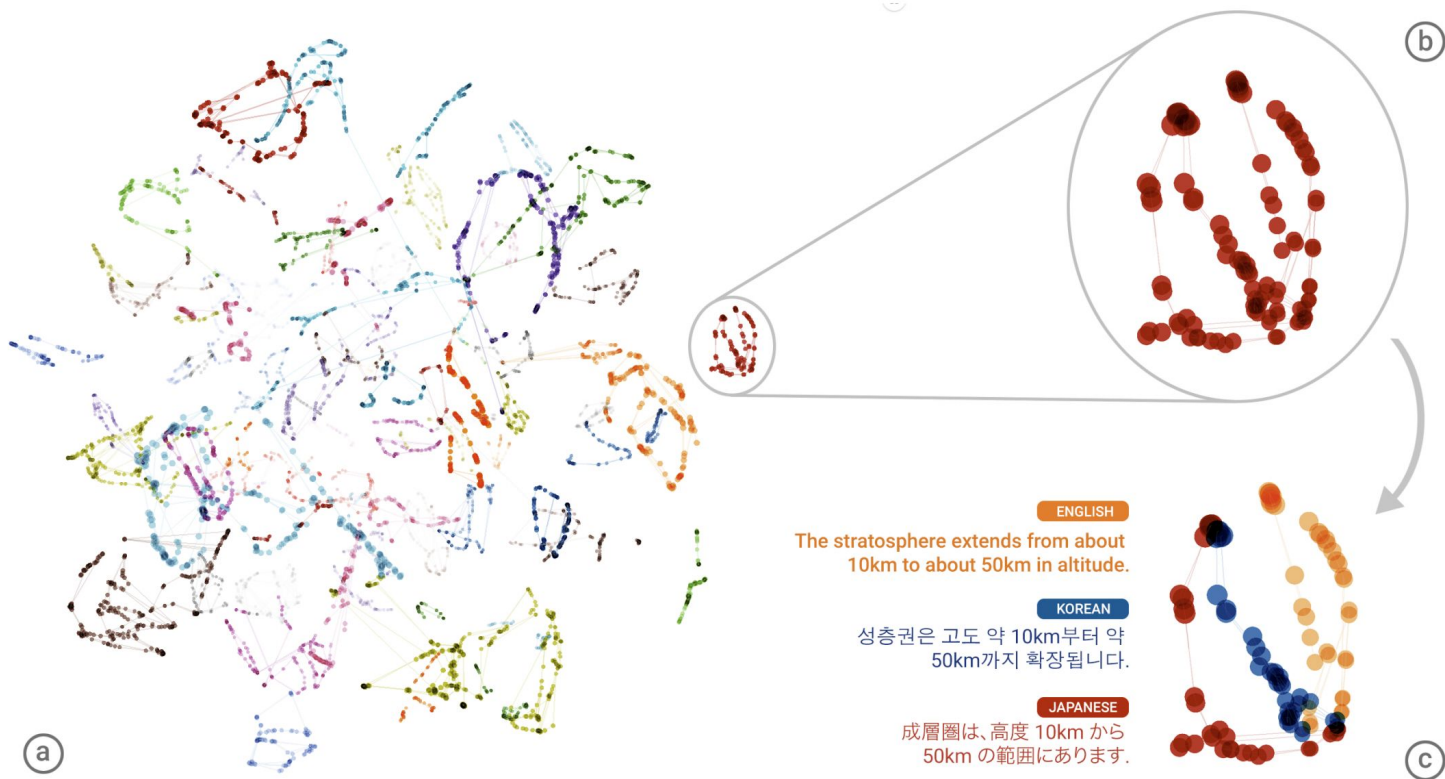
... where every other token is OOV



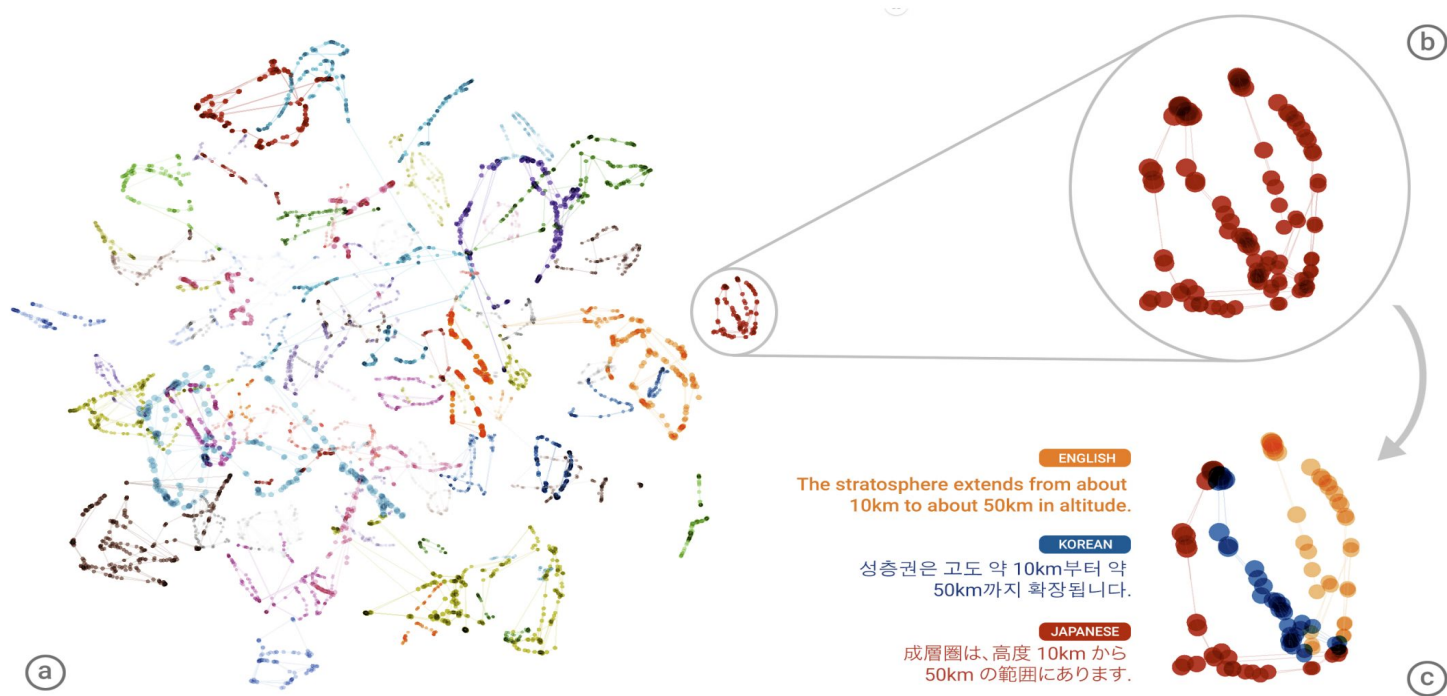
... and you'll learn lessons that will apply to the
>6700 human languages that are not like English



Research questions

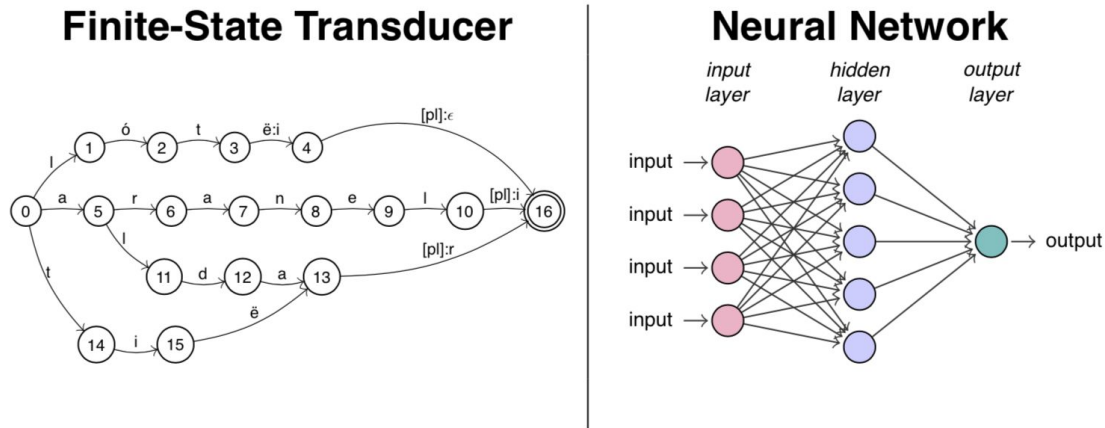


1) What embedding representations are most appropriate for languages with a high ratio of morphemes to words?



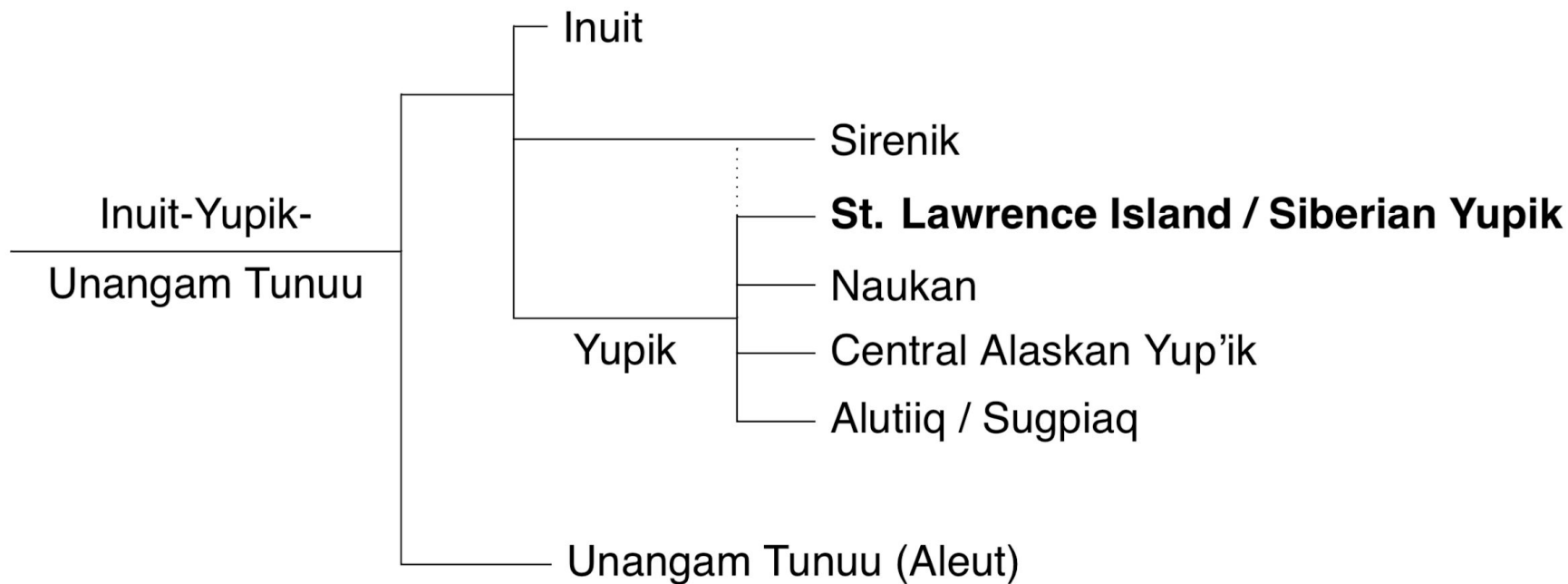
2) Can FST-based morphological analyzers effectively be used to bootstrap *more* effective neural LMs?

- Morphological analyzers may be implemented as a



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

3) How can common neural models be developed for closely related languages in order to maximize the utility of sparse digitized resources?

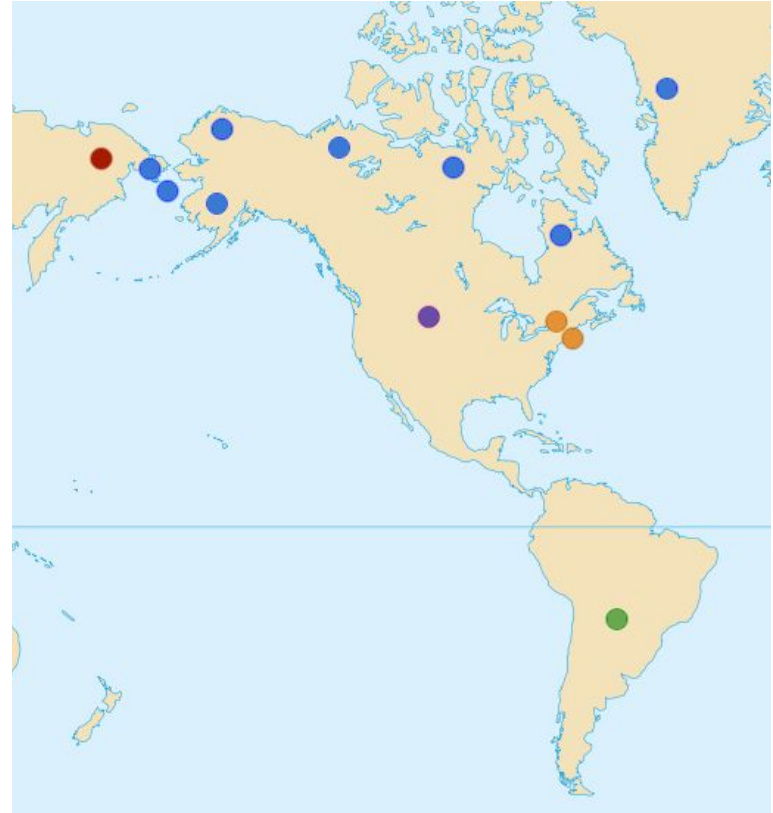


Overview

- Language selection and data collection
- Baseline systems
 - Language modelling
 - Machine translation
 - ASR
- Downstream applications
 - Spell-checkers & on-device morph-aware dictionaries
 - On-device predictive text
 - Audio transcription / search of audio archives
 - Interactive machine translation of monolingual polysynthetic texts

Language selection

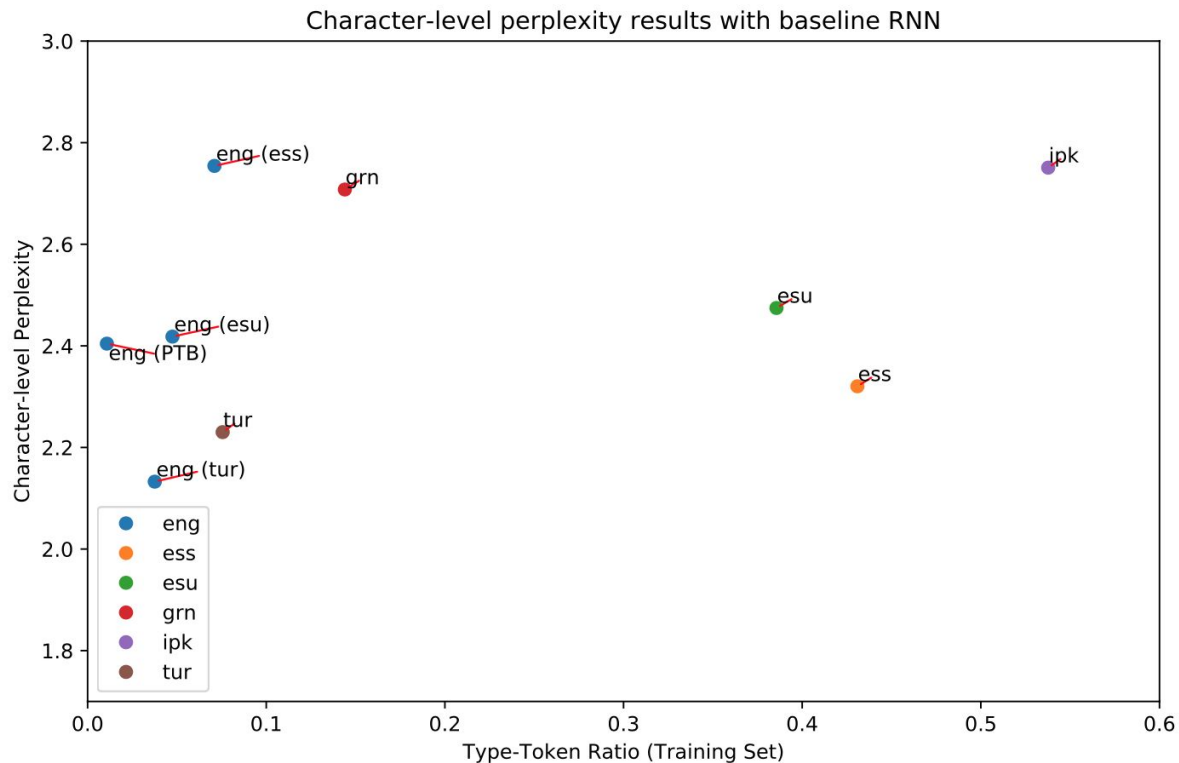
Saint Lawrence Island Yupik
Alaskan Yup'ik
North Slope Iñupiaq
Inuktitut
Kalaallisut
Chukchi
Guaraní
Seneca
Mohawk
Crow



Data collection

Language	Code	Monolingual	Parallel	FST	Audio
St. Lawrence Island Yupik	ess	24,456	8,002	✓	33h*
Alaskan Yup'ik	esu	—	45,254	—	—
North Slope Iñupiaq	esi	4,070	—	✓	—
Inuktitut	iku	1,552	1,300,159	(✓)	—
Kalaallisut	kal	5,949	—	✓	—
Chukchi	ckt	1,015	—	✓	1h30*
Guaraní	grn	—	30,078	✓	—
Seneca	see	—	—	—	12h
Mohawk	moh	—	—	—	—
Crow	cro	—	—	—	7h

Baseline system: RNN LM



Baseline system: MT

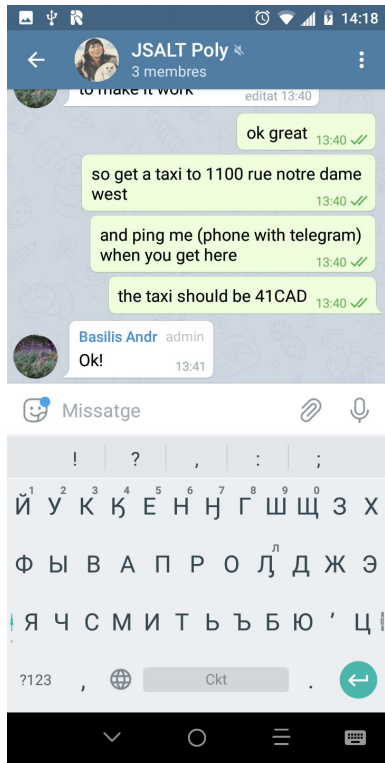
English → Inuktitut

BPE	NMT	SMT
500	16.7	13.1
1000	16.8	14.8
5000	16.9	15.2
15000	17.0	15.1
30000	16.7	15.2
60000	16.8	15.5

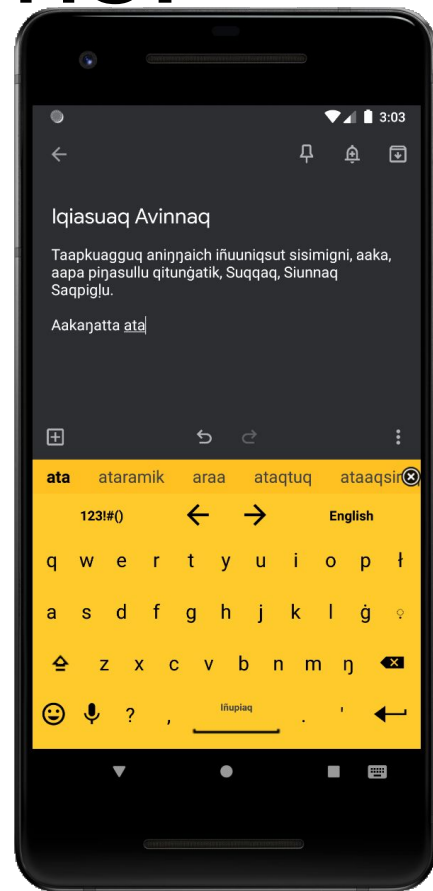
Baseline system: ASR

- DeepSpeech ASR baseline for Crow
- Forced aligner trained for SLI Yupik

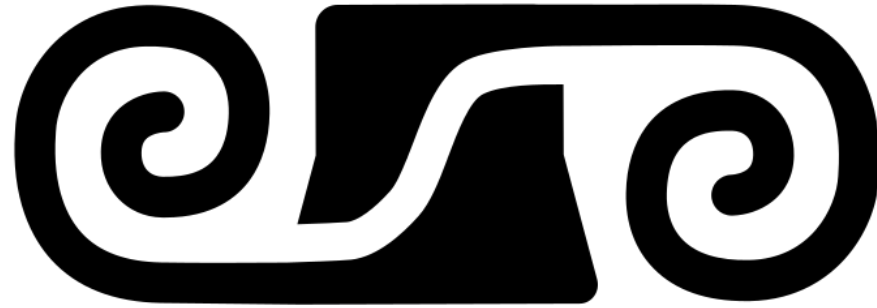
Downstream applications: On-device predictive text




Giellatekno



Downstream applications:
Spell-checkers



Giellatekno

Downstream applications: Interactive MT

The screenshot displays a web-based interface for Interactive Machine Translation (MT). It features a header with Inuktitut text, a main content area with English text on the left and Inuktitut text on the right, and a search bar at the bottom. The search bar contains the Inuktitut text "sivulliqaami | umiarjuakkut tikititauvattunut". Below the search bar are four buttons: "ITP", "SRC→", "DRAFT", and "TRANSLATED". The "TRANSLATED" button is highlighted in blue. A "visualization >>" link is located in the top right corner of the main content area.

samanit taqqinik iqqanaijaqsinnaat ilinniaqtitaullutik iqalunni.

staff was a valuable one for the staff and the department as a

tamanna ilinniaqtittiniq nunalimmiunit iqqanaijaqtinit piujummariulauqupuq iqqanaijaqtinut am-
ma pilirivivilimaamut.

ince the Eastern Arctic Sealift Program was transferred to
t from the Canadian Coast Guard has just recently ended.

visualization >>

sivulliqaami | umiarjuakkut tikititauvattunut

ITP SRC→ DRAFT TRANSLATED

Downstream applications: Audio transcription / search



Workshop goals

- What embedding representations and neural LM architectures are most effective for languages with a high ratio of morphemes per word?
- Can rule-based FSTs be used to bootstrap more robust neural models?
- To what extent can common neural models be developed for closely related languages in order to maximize the utility of sparse digitized resources?

NPLM Team Members

- **Lane Schwartz (University of Illinois)**
- **Francis Tyers (Indiana University)**
- Lori Levin (CMU)
- Christo Kirov (Google AI)
- Patrick Littell (NRC Canada)
- Jackie Lo (NRC Canada)
- Emily Prud'hommeaux (Boston College)
- **Hayley Park (Illinois)**
- **Kenneth Steimel (Indiana)**
- **Lonny Strunk (Washington)**
- **Coleman Hayley (JHU)**
- **Katherine Zhang (CMU)**
- Jeffrey Micher (CMU)
- Rebecca Knowles (JHU)
- Robbie Jimmerson (BC)
- Vasya Adriyanets (HSE)
- Han Liu (UC Boulder)

Let's learn lessons that will apply to the >6700
human languages that are not like English

