

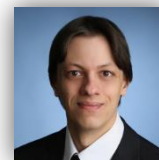
Ethics in Natural Language Processing – SS 2022



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Lecture 7 Low-Resource NLP

Dr. Thomas Arnold
Aniket Pramanik



Ubiquitous Knowledge Processing Lab
Technische Universität Darmstadt

Slides and material from Yulia Tsvetkov



Carnegie Mellon University
Language Technologies Institute

Syllabus (tentative)

<u>Nr.</u>	<u>Lecture</u>
01	Introduction, Foundations I
02	Foundations II
03	Bias I
04	Bias II
05	Incivility and Hate Speech I
06	NO LECTURE – Christi Himmelfahrt
07	Incivility and Hate Speech II
08	Low-Resource NLP
09	NO LECTURE - Fronleichnam
10	Privacy and Security I
11	Privacy and Security II
12	Language of Manipulation I
13	Language of Manipulation II

Learning Goals

After hearing this lecture, you should be able to...

- **Explain why traditional NLP techniques fail on other languages**
- **Explain syntactic and semantic ambiguities**
- **Discuss the advantages and disadvantages of statistical neural NLP methods in Low-Resource scenarios**

Low-Resource NLP: Introduction

Approaches to Low-Resource NLP

Low-Resource NLP: Lorelei Program

What does an NLP system need to “know”?



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Language consists of many levels of structure

Humans fluently integrate all of these in producing/understanding language

Ideally, so would a computer!

Sounds

SOUNDS

Th i a si e n

Words

WORDS

This is a simple sentence

Morphology

WORDS

MORPHOLOGY

This is a simple sentence

be
3sg
present

Parts of Speech

PART OF SPEECH

WORDS

MORPHOLOGY

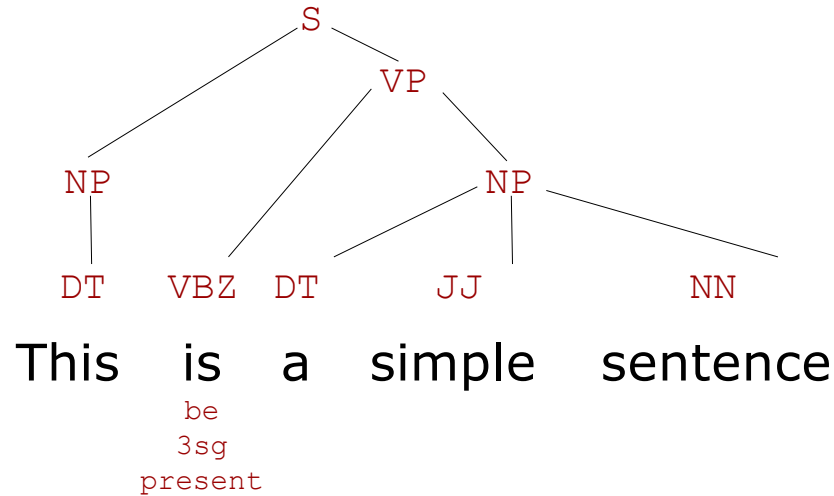
DT	VBZ	DT	JJ	NN
This	is	a	simple	sentence
	be			
	3sg			
	present			

SYNTAX

PART OF SPEECH

WORDS

MORPHOLOGY



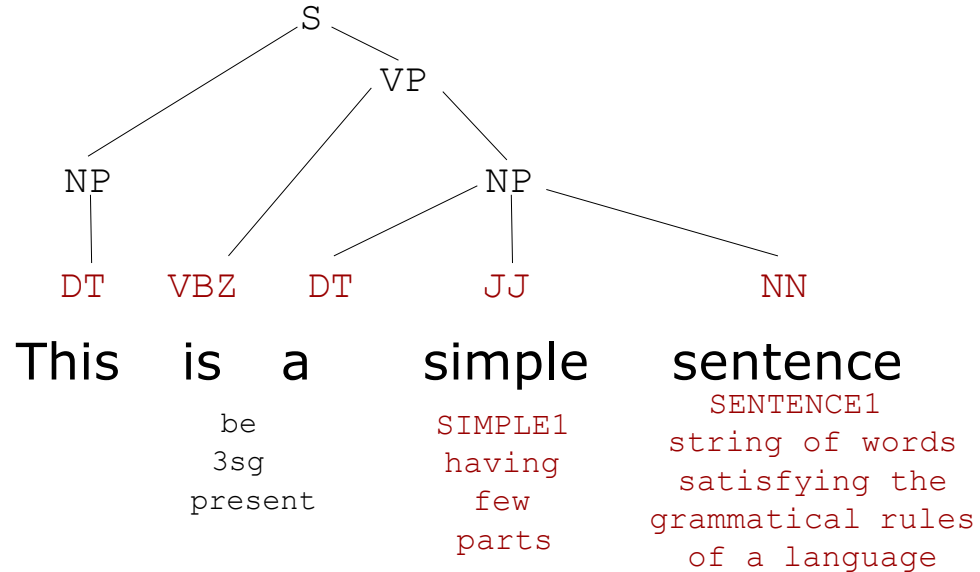
SYNTAX

PART OF SPEECH

WORDS

MORPHOLOGY

SEMANTICS



SYNTAX

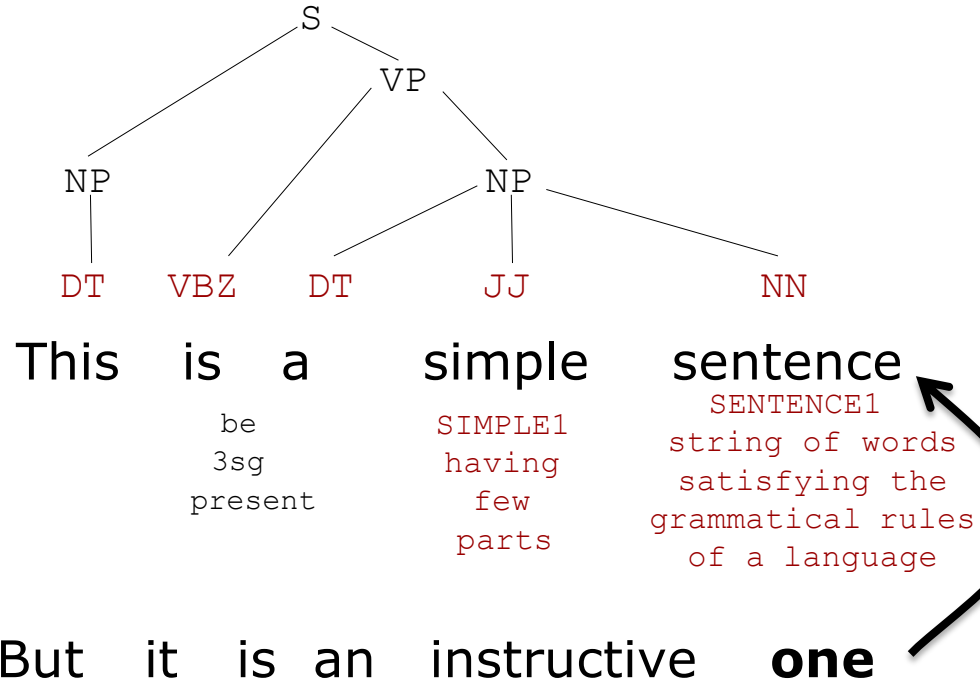
PART OF SPEECH

WORDS

MORPHOLOGY

SEMANTICS

DISCOURSE



- Core technologies

- Language modelling
- Part-of-speech tagging
- Syntactic parsing
- Named-entity recognition
- Coreference resolution
- Word sense disambiguation
- Semantic Role Labelling
- . . .

- Applications

- Machine Translation
- Information Retrieval
- Question Answering
- Dialogue Systems
- Information Extraction
- Summarization
- Sentiment Analysis
- . . .

Why NLP is hard?

1. Ambiguity at many levels:

- Word senses: **bank** (finance or river?)
- Part of speech: **chair** (noun or verb?)
- Syntactic structure: **I saw a man with a telescope**
- Quantifier scope: **Every child loves some movie**
- Multiple: **I saw her duck**

⇒ NLP algorithms model ambiguity, and choose the correct analysis in context

2. Linguistic diversity

Why NLP is hard?

1. Ambiguity at many levels:

- Word senses: **bank** (finance or river?)
- Part of speech: **chair** (noun or verb?)
- Syntactic structure: **I saw a man with a telescope**
- Quantifier scope: **Every child loves some movie**
- Multiple: **I saw her duck**

⇒ NLP algorithms model ambiguity, and choose the correct analysis in context

2. Linguistic diversity

6 – 7k World Languages



Linguistic Diversity: Words

这是一个简单的句子

WORDS

This is a simple sentence

זה משפט פשוט

Linguistic Diversity: Hebrew Words

in tea
her daughter

בתה

- most of the vowels unspecified

Linguistic Diversity: Words

התבשו

and her saturday

ו +תבש+ה

and that in tea

ו +ש+ב+הת

and that her daughter

ו +ש+תב+ה

- most of the vowels unspecified
- particles, prepositions, the definite article, conjunctions attach to the words which follow them
- tokenization is highly ambiguous

Linguistic Diversity: Morphology

WORDS

This is a simple sentence

MORPHOLOGY

be
3sg
present

Linguistic Diversity: Quechua Morphology

Much'anayanayakapushasqakupuniñataqsunamá
Much'a -na -naya -ka -pu -sha -sqa -ku -puni -ña -taq -suna -má

"So they really always have been kissing each other then"

Much'a to kiss

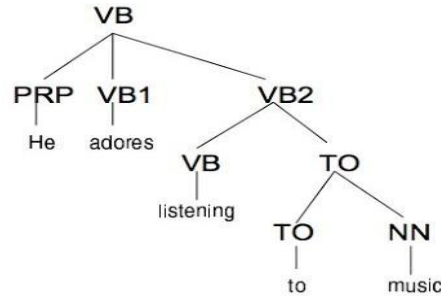
- na expresses obligation, lost in translation
- naya expresses desire
- ka diminutive
- pu reflexive (kiss *eachother*)
- sha progressive (kiss*ing*)
- sqa declaring something the speaker has not personally witnessed
- ku 3rd person plural (they kiss)
- puni definitive (really*)
- ña always
- taq statement of contrast (...then)
- suna expressing uncertainty (So...)
- má expressing that the speaker is surprised

Linguistic Diversity: Russian Morphology

	Singular+neut	Plural+neut	
Nominative	предложение	предложения	sentence (s)
Genitive	предложения	предложений	(of) sentence (s)
Dative	предложению	предложениям	(to) sentence (s)
Accusative	предложение	предложения	sentence (s)
Instrumental	предложением	предложениями	(by) sentence (s)
Prepositional	предложении	предложениях	(in/at) sentence (s)

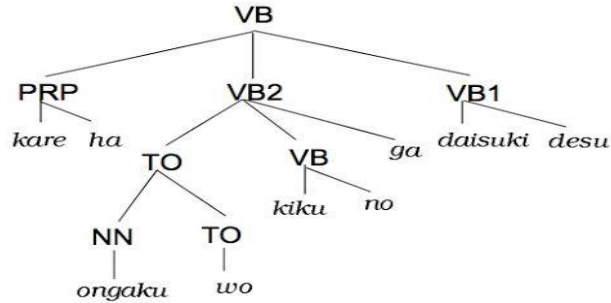
Linguistic Diversity: Japanese Syntax

SVO



he adores listening to music

SOV

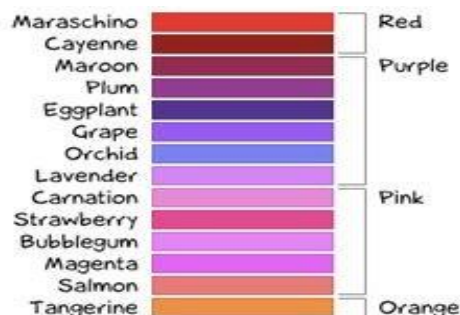


かれはおんがくをきくのがだいすきです
kare ha ongaku wo kiku no ga daisuki desu

he adores listening to music

(Yamada & Knight '02)

Linguistic Diversity: Semantics



Russian has relatively few names for colors; Japanese has hundreds

Multiword expressions, e.g. **it's raining cats and dogs** or **wake up** and metaphors, e.g. **Love is a journey** are very different across languages

Why NLP is hard?

1. Ambiguity at many levels:

- Word senses: **bank** (finance or river?)
- Part of speech: **chair** (noun or verb?)
- Syntactic structure: **I saw a man with a telescope**
- Quantifier scope: **Every child loves some movie**
- Multiple: **I saw her duck**

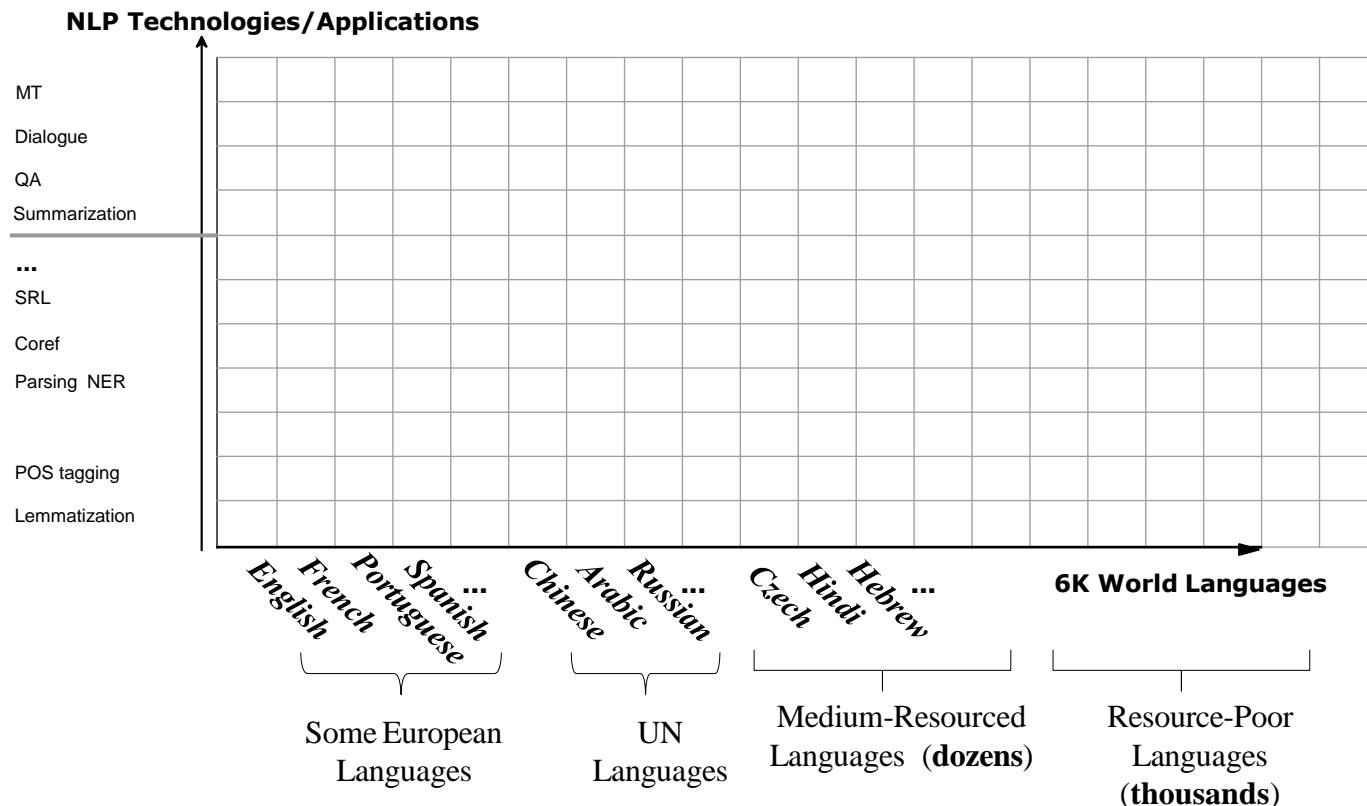
⇒ NLP algorithms model ambiguity, and choose the correct analysis in context

2. Linguistic diversity

- 6–7K languages in the world
- Languages diverge across all levels of linguistic structure
 - ⇒ **no generic solution for a particular NLP task**
- **Most of the languages do not have sufficient resources to build statistical NLP models**

Low-resource languages— languages lacking large monolingual or parallel corpora and/or manually crafted linguistic resources sufficient for building statistical NLP applications

What NLP Technologies are Resource-Rich?



Performance of Resource-Rich vs. Resource-Poor NLP

Machine Translation

Parallel corpus

Nenhum deles reparou na janela , através da qual teria podido ver uma enorme coruja amarelada , esvoaçando em grande alvoroço .

assim , não viu as corujas descendo rapidamente em plena luz do dia , apesar de todos os transeuntes apontarem estarecidos e de boca aberta enquanto coruja após coruja lhes passavam A grande velocidade sobre as cabeças .

Queira enviar-nos A sua coruja até dia 31 de Julho , sem falta .

- O que é que quer dizer esperarem A minha coruja ?

Hagrid Hagrid enrolou A nota , deu-a à coruja que A agarrou com O bico e , dirigindo-se à porta , soltou A ave no meio da tempestade .

O próprio Hagrid adormecera no sofá totalmente destruído e , bicando no vidro da janela , estava uma coruja que segurava um jornal .

A coruja entrou e depôs O jornal em cima de Hagrid

None of them noticed a large , tawny owl flutter past the window .

He didn ' t see the owls swoop ing past in broad daylight , though people down in the street did ; They pointed and gazed open-mouthed as owl after owl sped overhead .

We await your owl by no later than July 31 .

after a few minutes He stammered , " what does it mean , They await My owl ?

Hagrid Hagrid rolled up the note , gave it to the owl , which clamped it in its beak , went to the door , and threw the owl out into the Storm .

the hut was full of sunlight , the Storm was over , Hagrid himself was asleep on the collapsed sofa , and there was an owl rapping its claw on the window , a newspaper held in its beak .

the owl swooped in and dropped the newspaper on top of Hagrid , who didn ' t

Resource-rich: millions of parallel sentences

Resource-poor: few thousands of parallel sentences

Performance of Resource-Rich vs. Resource-Poor NLP Machine Translation

English → French

Translate Turn off instant translation

Russian English French Detect language

English Spanish French Translate

You will just have to find a way of getting over it. x

Vous devrez trouver un moyen de le surmonter.

52/5000 Suggest an edit

French → English

Translate Turn off instant translation

Russian English French Detect language

English Spanish French Translate

Vous devrez trouver un moyen de le surmonter. x

You will have to find a way to overcome it.

45/5000 Suggest an edit

Did you mean: Vous devez trouver un moyen de le surmonter.

Performance of Resource-Rich vs. Resource-Poor NLP Machine Translation

English → Swahili

Translate Turn off instant translation

Russian English French Detect language

English Swahili French Translate

You will just have to find a way of getting over it.

Utakuwa tu kupata njia ya kupata juu yake.

Suggest an edit

53/5000

Swahili → English

Translate Turn off instant translation

Swahili English French Detect language

English Swahili French Translate

Utakuwa tu kupata njia ya kupata juu yake.

You will just find the way to get on it.

Suggest an edit

42/5000

Performance of Resource-Rich vs. Resource-Poor NLP Machine Translation

English → Hindi → English

Hindi English Yoruba Detect language

English Yoruba Hindi Translate

आपको इसे खत्म करने का एक तरीका मिलना होगा।

You have to find a way to eliminate it.

42/5000

Suggest an edit

English → Telugu → English

Uzbek English Telugu Detect language

English Uzbek Telugu Translate

మీరు దాని పైకి రావడానికి ఒక మార్గాన్ని కనుగొనవలసి ఉంటుంది.

You have to find a way to get it up.

59/5000

Suggest an edit

English → Uzbek → English

Pashto English Uzbek Detect language

English Uzbek Yoruba Translate

Buning ustiga faqatgina bir usulni topish kerak.

On top of that, you just have to find a way out.

48/5000

Suggest an edit

Performance of Resource-Rich vs. Resource-Poor NLP Machine Translation

English → Swahili

Swahili English Telugu Detect language ▾

English Swahili Telugu ▾ Translate

The summer school is meant to be an introduction to the state-of-the-art research in the speech and language technology area for graduate and undergraduate students.

Shule ya majira ya joto ina maana ya kuanzishwa kwa utafiti wa hali ya sanaa katika eneo la teknolojia na lugha ya wanafunzi kwa wanafunzi wahitimu na wahitimu.

166/5000

Suggest an edit

Swahili → English

Swahili English Telugu Detect language ▾

English Swahili Telugu ▾ Translate

Shule ya majira ya joto ina maana ya kuanzishwa kwa utafiti wa hali ya sanaa katika eneo la teknolojia na lugha ya wanafunzi kwa wanafunzi wahitimu na wahitimu.

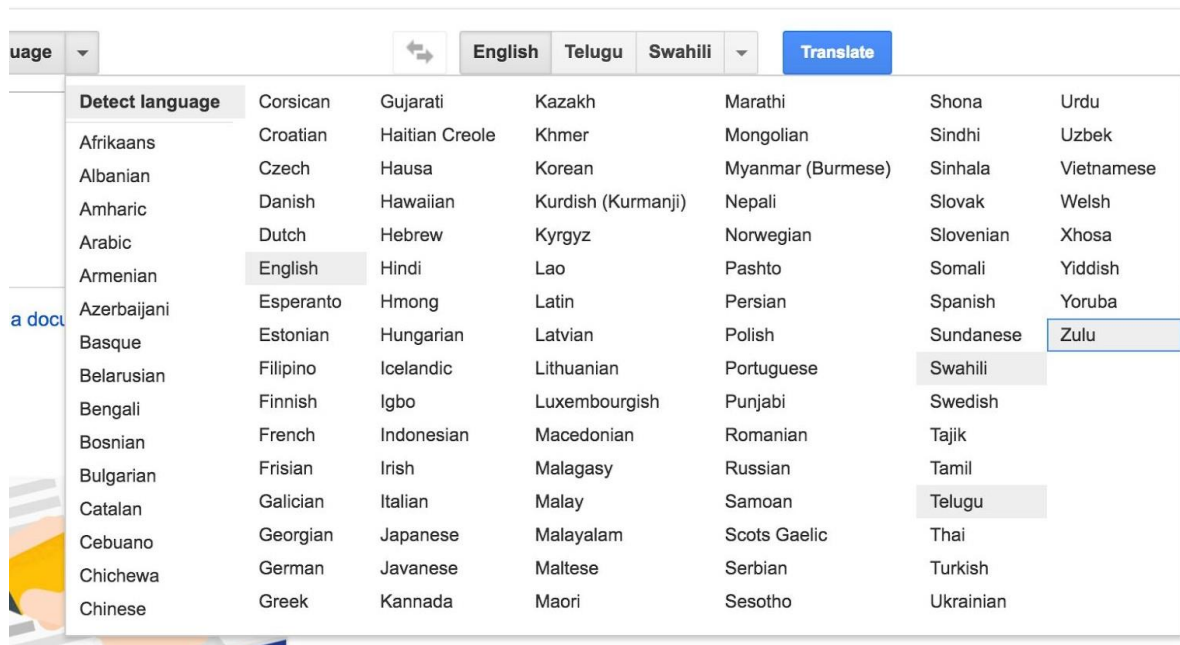
Summer school means the establishment of a state-of-the-art arts research technology and pupil language for graduate students and graduates.

160/5000

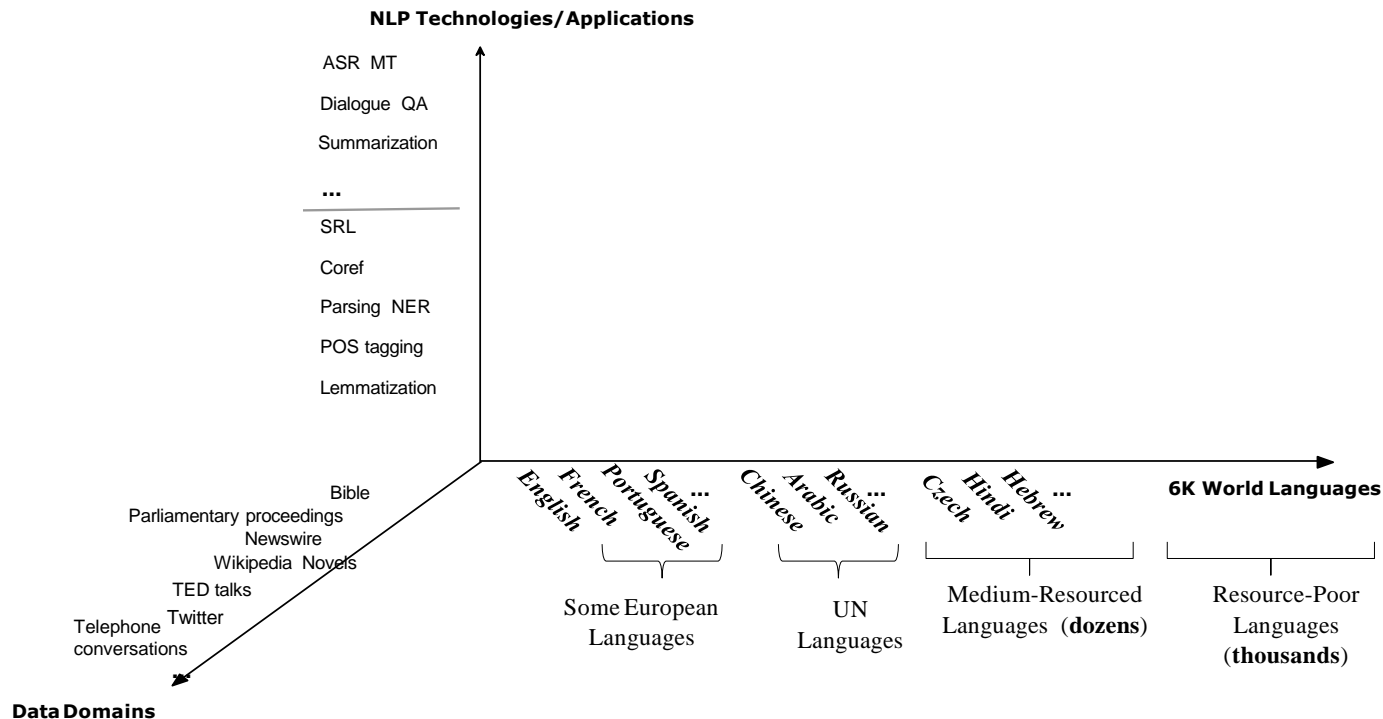
Suggest an edit

Performance of Resource-Rich vs. Resource-Poor NLP Machine Translation

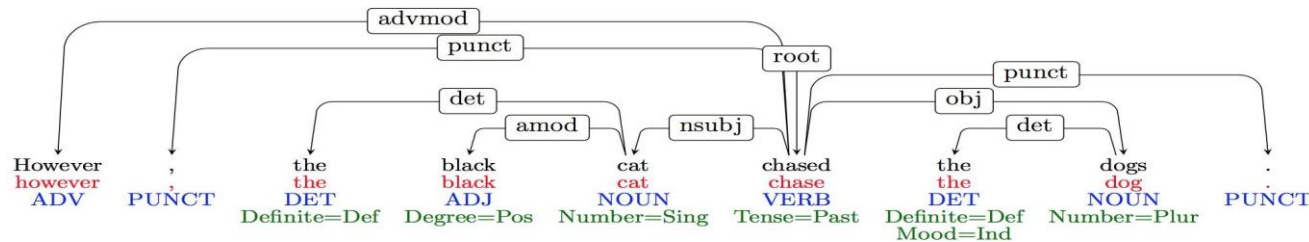
About 130 out of 6K languages (as of 2022)



Low-Resource NLP is Not Only About Multilinguality



Parsing models trained using Wall Street Journals treebanks do not work for spoken language domain



Spoken language is riddled with verbal disfluencies that interrupt the flow of speech, including long pauses, repeated words or phrases, restarts, and revisions of content:

Um, the black the black cat ch- chased the dogs.

Low-Resource NLP is Not Only About Multilinguality

Twitter processing is hard...



 Follow

Like serious dis flu nor dey wan go oooo.... Sick



Venus
@christinedarvin

 Follow

@_rkpntnte hindi ko alam babe eh, absent ako
kanina I'm sick rn hahaha 🤔👏



Donald J. Trump 
@realDonaldTrump

 Follow

Despite the constant negative press covfefe

Low-Resource NLP is Not Only About Multilinguality

- Much of the world knowledge is not in text, corpora contain what people said, but not what they meant, or how they understood things, or what they did in response to the language

This is milk



?



Outline

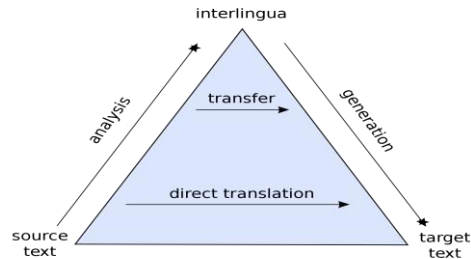
Low-Resource NLP: Introduction

Approaches to Low-Resource NLP

Low-Resource NLP: Lorelei Program

Paradigm Shifts in NLP

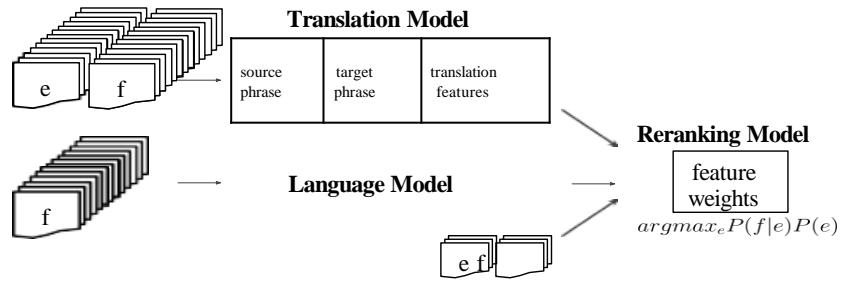
Logic-based/Rule-based NLP



~90s



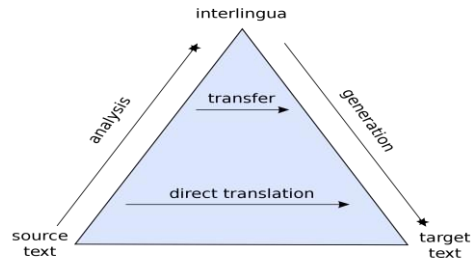
Statistical NLP



Rule-based models: high precision but very low recall

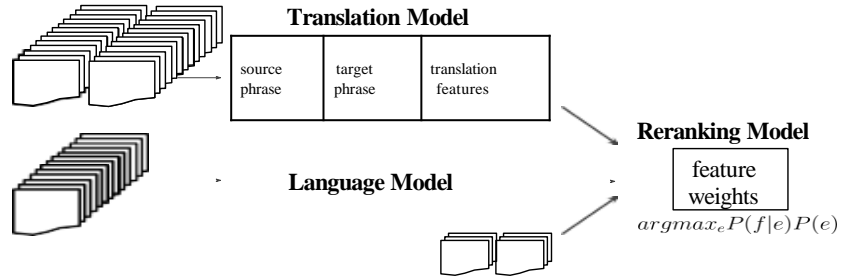
Paradigm Shifts in NLP

Logic-based/Rule-based NLP



* In resource-rich settings

Statistical NLP



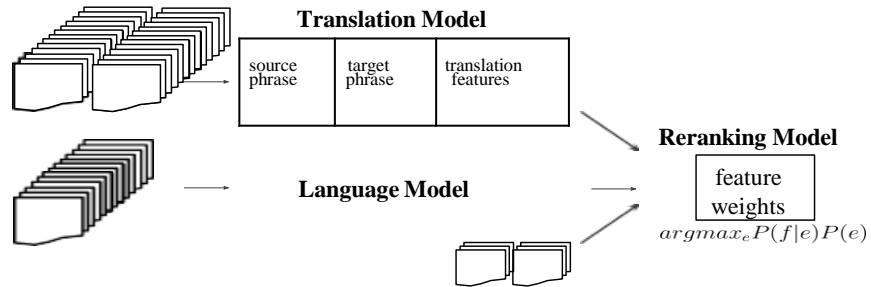
Statistical models: robust in the face of real-world data

Better performance

Less engineering of hand-crafted rules/knowledge

Paradigm Shifts in NLP

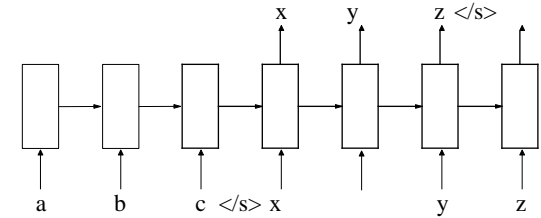
Statistical NLP



~mid 2010s

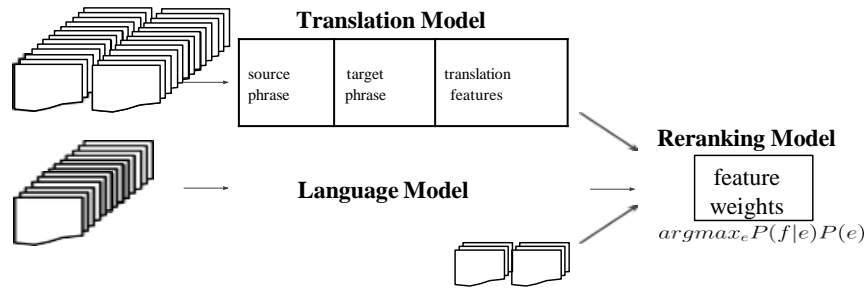


Statistical Neural NLP



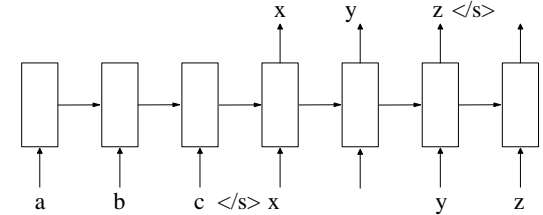
Paradigm Shifts in NLP

Statistical NLP



* In resource-rich settings

Statistical Neural NLP



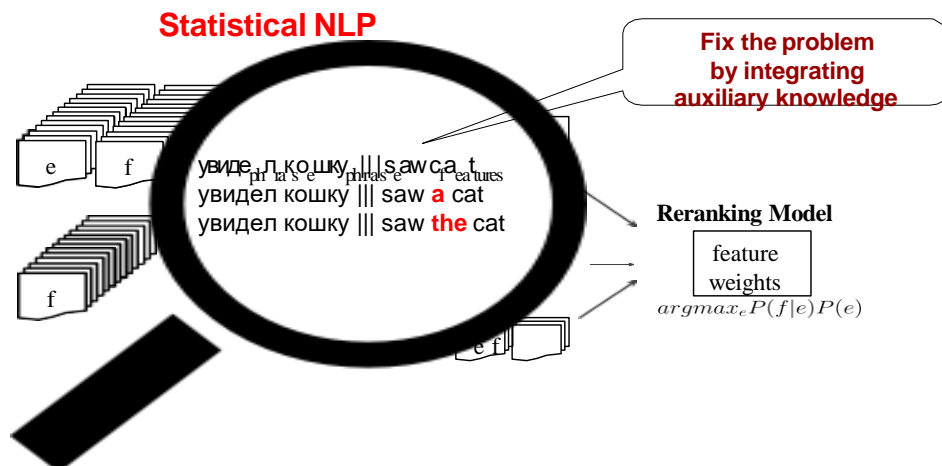
Robustness in the face of real-world data

Better performance

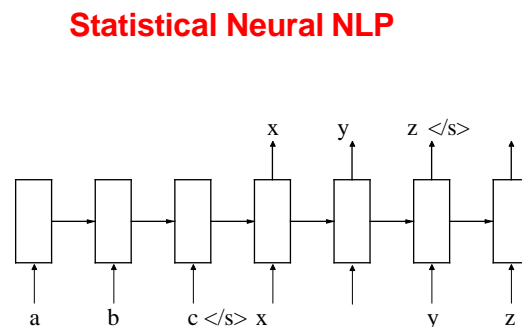
Less engineering of hand-crafted rules/knowledge

Building Blocks in Conventional Statistical NLP Models

Words, phrases \Rightarrow easier analysis, easier adaptation

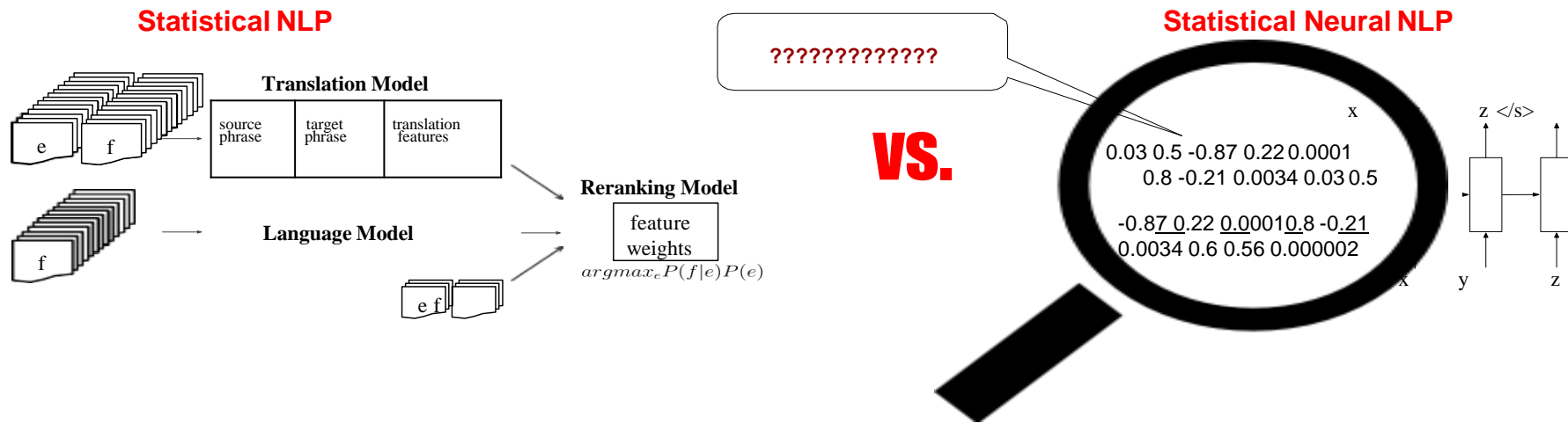


VS.



Building Blocks in Neural NLP Models

Vectors, matrices \Rightarrow not clear yet how to interpret



How to interpret continuous representations?

How to integrate auxiliary knowledge into neural network architectures?

- State-of-the-art NLP models require large amounts of training data and/or sophisticated language-specific engineering
- Large amounts of training data are unavailable for most languages
 - extreme case: languages that don't have a written form, e.g. Shanghainese spoken by 14 million people
 - or languages that just don't have online presence, e.g. Chichewa, a Bantu language spoken by 12 million people
- Language-specific engineering is expensive, requires linguistically trained speakers of the language

Unsupervised Learning

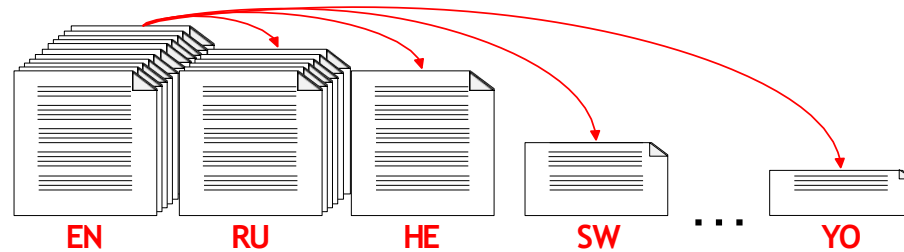
Unsupervised feature induction: Brown clustering, Word vectors

Unsupervised POS tagging

Unsupervised dependency parsing

. . .

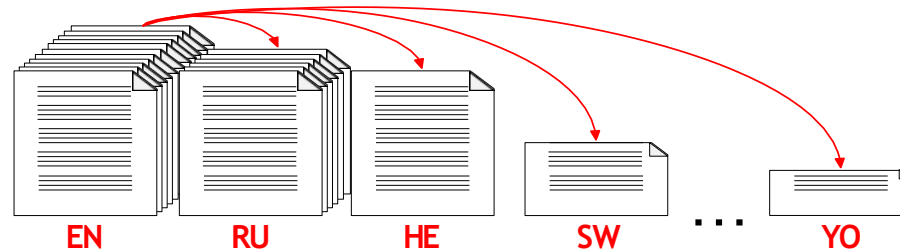
Transfer Learning



Cross-lingual transfer learning – transfer of resources and models from resource-rich source to resource-poor target languages

- Transfer of annotations (e.g., POS tags, syntactic or semantic features) via cross-lingual bridges (e.g., word or phrase alignments)
- Transfer of models – train a model in a resource-rich language and apply it in a resource-poor language

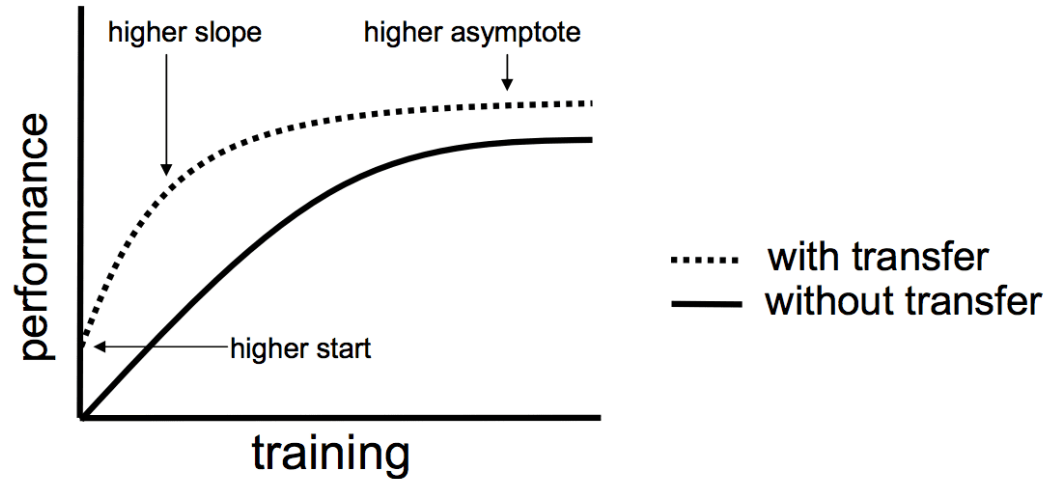
Transfer Learning



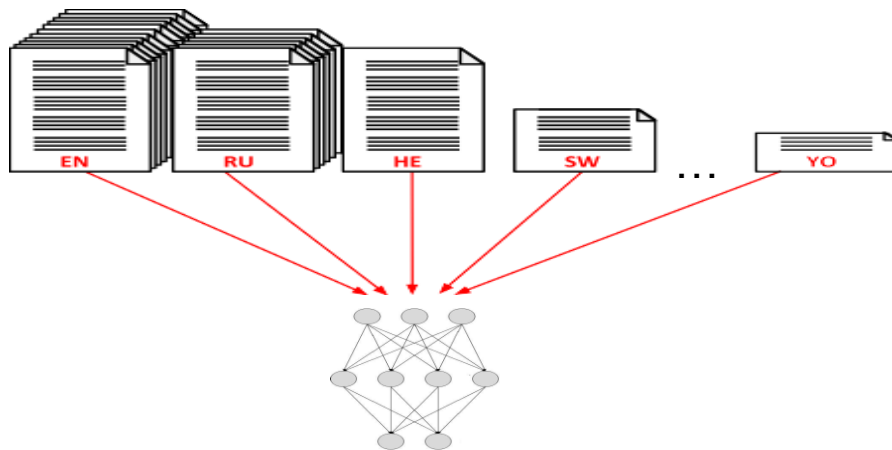
Zero-shot learning – train a model in one domains and assume it generalizes more or less out-of-the-box in a low-resource domain

One-shot learning – train a model in one domain and use only few examples from a low-resource domain to adapt it

Transfer Learning: 3 Advantages



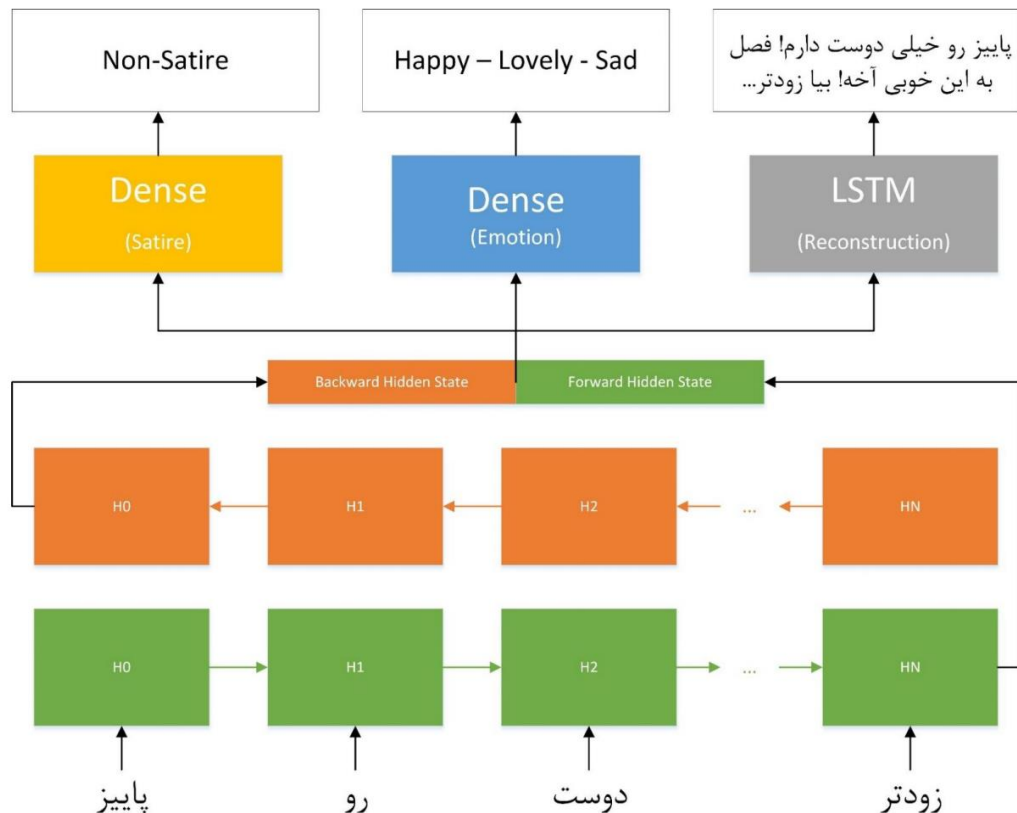
Joint Multilingual or “Polyglot” Learning



Joint resource-rich and resource-poor learning using a language-universal representation.

- Convert data in all languages to a shared representation
- Train a single model on a mix of datasets in all languages, to enable parameter sharing where possible

Case Study: Satire Detection



Satire Dataset: 2K tweets
Emotion Dataset: 300K
Reconstruction Dataset: >200M

Satire Performance (F1)

Single Task: 0.55
Multi-Task: 0.68

Another Approach: Task Modeling

Basic Idea: Model your problem in an easy-to-acquire-label setting

Example: Sentiment / Emotion analysis

Problem: People do not annotate emotions in text very often...

But: People use emoji all the time!

Solution: Approximate emotion prediction by emoji prediction!

Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm:

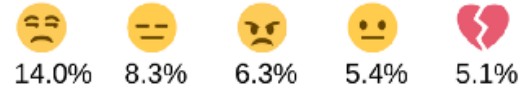
<https://arxiv.org/pdf/1708.00524.pdf>

Emoji: Proxy for Emotions

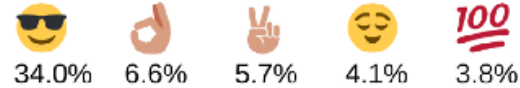
I love mom's cooking



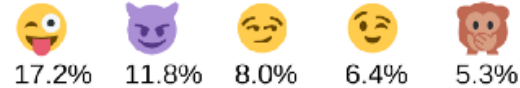
I love how you never reply back..



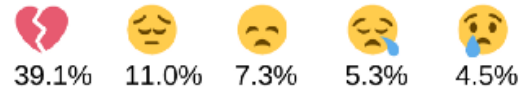
I love cruising with my homies



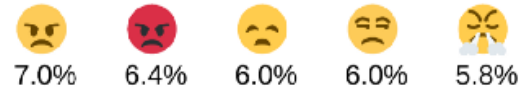
I love messing with yo mind!!



I love you and now you're just gone..



This is shit



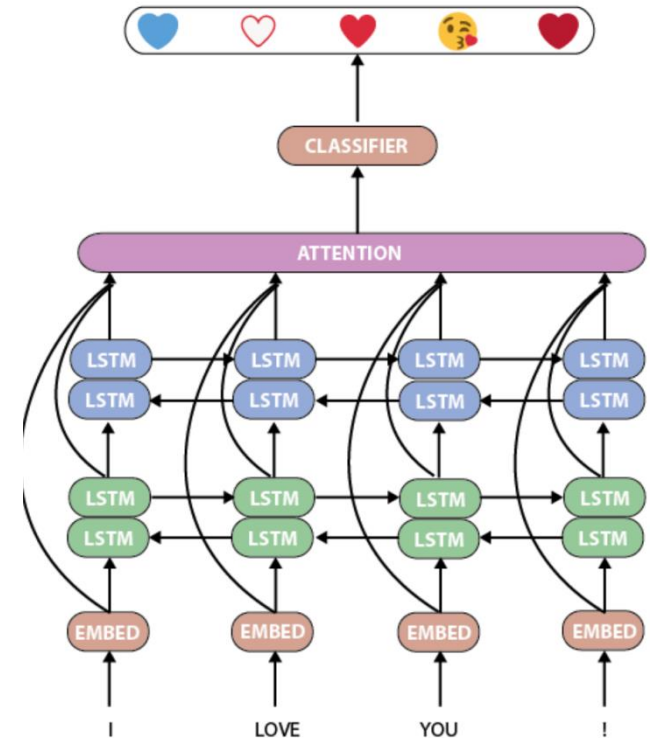
This is the shit



Emoji Classifier

DeepMoji Model

- Predict Emoji
- Map Emoji to Emotions



Source: http://sharif.edu/~kharrazi/data_talks/slides/data-talks-sabeti-98-7-24-slides.pdf

Emoji Prediction

پاییز رو خیلی دوست دارم! فصل
به این خوبی آخه! بیا زودتر...



بی حس

5.46%



ناراحت

9.28%



عشق

34.41%



خوشحال

39.04%

Summary

- Low-Resource NLP is hard!
- Ambiguity (word senses, part-of-speech, syntactic structure...)
- Linguistic diversity at all levels of language structure
 - Tokenization, morphology, part-of-speech, syntax, semantics, discourse...
- Paradigm shifts in NLP
 - Rule-based NLP: high precision, low recall
 - Statistical NLP: Needs more data
 - Neural NLP: Needs MORE more data!
- Most promising approach: Transfer Learning

Low-Resource NLP: Introduction

Approaches to Low-Resource NLP

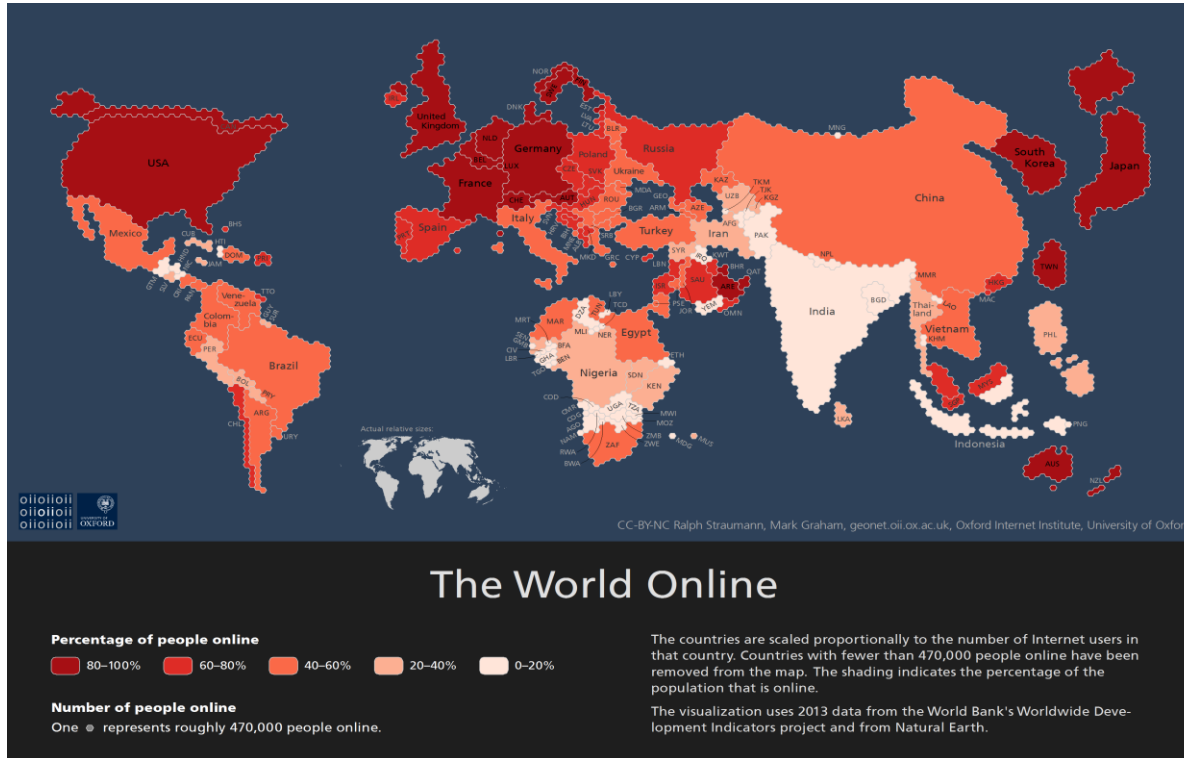
Low-Resource NLP: Lorelei Program

Why Care About Low-Resource NLP?

1. Commercial value
2. Social-good reasons

Why Care About Low-Resource NLP?

Commercial value



Why Care About Low-Resource NLP?

Social good reasons

Translation systems

Speech interfaces

Dialogue systems

Educational applications

Emergency response applications

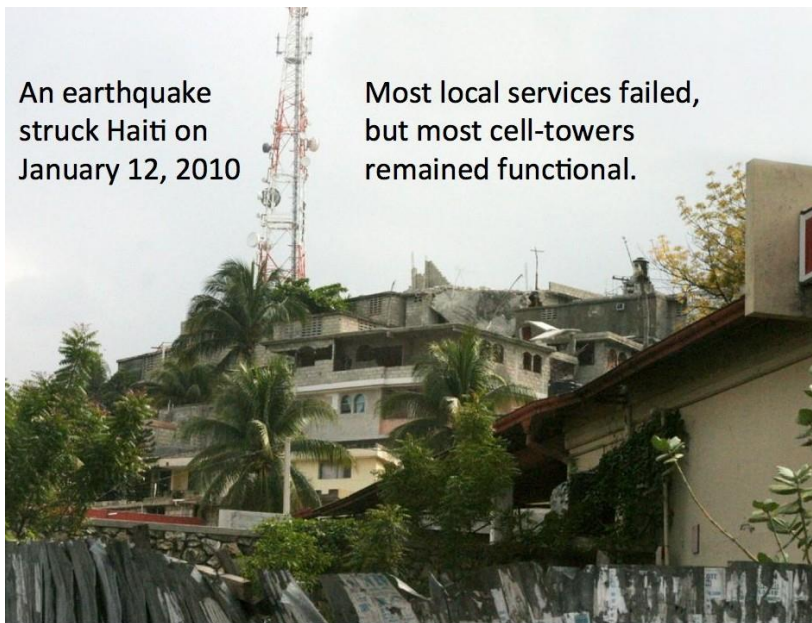
Monitoring democratic processes

Why Care About Low-Resource NLP?

Social good reasons

Social good reasons: Emergency response

About 3 million people were affected by the quake



Why Care About Low-Resource NLP?

Social good reasons: Emergency response

Messages start streaming in

- Fanmi mwen nan
Kafou, 24 Cote Plage,
41A bezwen manje
ak dlo
- Moun kwense nan
Sakre Kè nan
Pòtoprens
- Ti ekipman Lopital
General genyen yo
paka minm fè 24 è
- Fanm gen tranche
pou fè yon pitit nan
Delmas 31

iDiBON

* Slide by Rob Munro <http://web.stanford.edu/class/cs124>

Why Care About Low-Resource NLP?

Social good reasons: Emergency response

Messages start streaming in

- Fanmi mwen nan Kafou, 24 Cote Plage, 41A bezwen manje ak dlo
- Moun kwense nan Sakre Kè nan Pòtoprens
- Ti ekipman Lopital General genyen yo paka minm fè 24 è
- Fanm gen tranche pou fè yon pitit nan Delmas 31
- My family in Carrefour, 24 Cote Plage, 41A needs food and water
- People trapped in Sacred Heart Church, PauP
- General Hospital has less than 24 hrs. supplies
- Undergoing children delivery Delmas 31

iDIBON

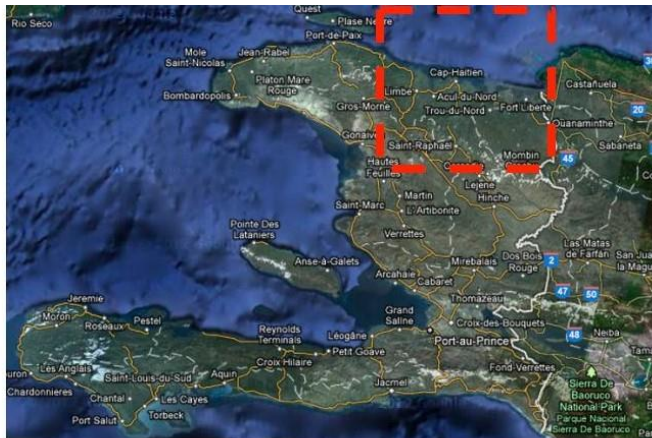
* Slide by Rob Munro <http://web.stanford.edu/class/cs124>

Why Care About Low-Resource NLP?

Social good reasons: Emergency response

Lopital Sacre-Coeur ki nan vil Okap, pre pou li resevwa moun malad e lap mande pou moun ki malad yo ale la.

“Sacre-Coeur Hospital which located in this village of **Okap** is ready to receive those who are injured. Therefore, we are asking those who are sick to report to that hospital.”



idibon

* Slide by Rob Munro <http://web.stanford.edu/class/cs124>

Why Care About Low-Resource NLP?

Identifying outbreaks of diseases



**Language
Detection**



Keyword Filter
“flu”, “sick”

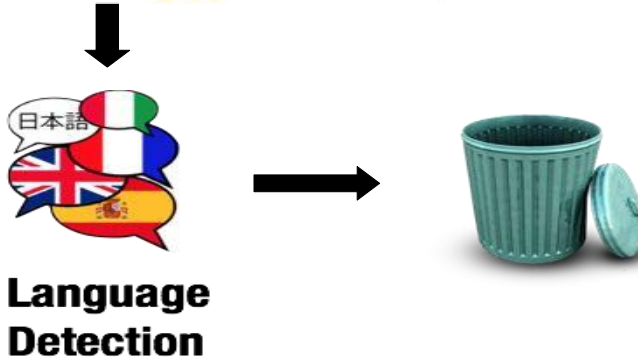


Analytics

Which symptoms?
Are they hungover?

Why Care About Low-Resource NLP?

Identifying outbreaks of diseases



Government Investment in Languages

- Language Technologies mostly developed for High Resource Languages
 - English, Spanish, German, Arabic, Mandarin
- What about the other 6995 languages?
 - Maybe 30 have good resources (ASR, Treebanks, Parsers)
- What about those around 300-1000?
 - > 1 Millions speakers, writing systems...
- If no immediate commercial value no support happens

US Government LT Investment

- DARPA (Defense Advance Research Projects Agency)
 - Invested in MT from 1940s
 - Invested in ASR from 1970s
 - Invested in Dialog systems from 1990s
 - Invested in Speech Translation from 1990s
- Case study Lorelei (2016-2021)

The Scenario

- Disaster happens! (e.g. earthquake)
- Area effected doesn't use major language
- Communication is in local language
 - News, TV/Radio, Social Media
- What is going on?
 - Where should you provide support
 - Who is affected
 - How many people need help
 - What is the urgency

Lorelei Incident

- Disaster happens! (e.g. earthquake)
- Communication is in local language
 - News, TV/Radio, Social Media
- Provide
 - Machine Translation
 - Named Entity Recognition
 - Situation Frames (11 types) plus location, status, urgency, “gravity”
 - evac, food, infra, med, search...

Lorelei Incident

- Disaster happens! (e.g. earthquake)
- Communication is in local language
 - News, TV/Radio, Social Media
- Provide
 - Machine Translation
 - Named Entity Recognition
 - Situation Frames (11 types) plus location, status, urgency, “gravity”
 - evac, food, infra, med, search...
- Do this in
 - 24 hours
 - 7 days
- You are told the language at hour 0

Lorelei Evaluation Exercises

- May 2016: Dry Run (Mandarin)
- July 2016: Uighur (Turkic Language spoken in Western China)
- July 2017: Tigrinya and Oromo (spoken in Eritrea and Ethiopia)
- July 2018: Kinyarwanda and Sinhala (spoken in Uganda and Sri Lanka)
- July 2019: ??? and ???

- Perform in pronunciation space
 - Not words, morphemes or character space
- Cross Lingual Transfer
 - If word3(L1) co-occurs with word1(L1), word2(L1)
 - And word3(L2) co-occurs with word1(L2), word2(L2)
 - And $\text{trans}(\text{word1}(\text{L1})) = \text{word1}(\text{L2})$ and $\text{trans}(\text{word2}(\text{L1})) = \text{word2}(\text{L2})$
 - Maybe $\text{trans}(\text{word3}(\text{L1})) = \text{word3}(\text{L2})$?

- Use available resources
 - Religious Texts (Bible, Quran, Unix Manuals...)
 - Wikipedia
 - Native Informant
- Global Linguistic Knowledge
 - High morphology language more likely to have free word order
 - Close language borrowing

- Techniques for low resource languages
 - Translation, interpretation, sentiment
 - Both particular languages, and general techniques
- Machine Learning
 - Better use of limited data
 - Use transfer learning
 - Not naive just end-to-end
 - Using large mono-lingual dataset to improve models
 - Using structure to make learning easier

In Two Weeks: Privacy & Security I