



Incorporating Dialectal Variability for Socially Equitable Language Identification

David Jurgens, Yulia Tsvetkov, and Dan Jurafsky



McNamee, P., “Language identification: *a solved problem* suitable for undergraduate instruction” *Journal of Computing Sciences in Colleges* 20(3) 2005.

McNamee, P., “Language identification: *a solved problem* suitable for undergraduate instruction” *Journal of Computing Sciences in Colleges* 20(3) 2005.

“This paper describes [...] how even the most simple of these methods **using data obtained from the World Wide Web achieve accuracy approaching 100%** on a test suite comprised of ten European languages”

Whose language are we identifying?

Whose language are we identifying?



The Royal Family ✓

@RoyalFamily

Follow



Taking place this week on the river Thames is 'Swan Upping' – the annual census of the swan population on the Thames.

Whose language are we identifying?



The Royal Family ✓
@RoyalFamily

Follow

Taking place this week on the river Thames is 'Swan Upping' – the annual census of the swan population on the Thames.



da'Rah-zingSun
@TIME7SS

Follow

[@kinguilfoyle](#) prblm I hve wit ur reportng is its 2 literal, evry1 knos pple tlk diffrent evrywhere, u kno wut she means jus like we do!

Whose language are we identifying?



The Royal Family ✓
@RoyalFamily

Follow

Taking place this week on the river Thames is 'Swan Upping' – the annual census of the swan population on the Thames.



da'Rah-zingSun
@TIME7SS

Follow

@kinguilfoyle prblm I hve wit ur reportng is its 2 literal, evry1 knos pple tlk diffnt evrywhere, u kno wut she means jus like we do!



Mooktar
@bossmukky

Follow

"@Ecstatic_Mi: @bossmukky Ebi like say I wan dey sick sef wlh 'Flu' my whole body dey weak"uw gee...



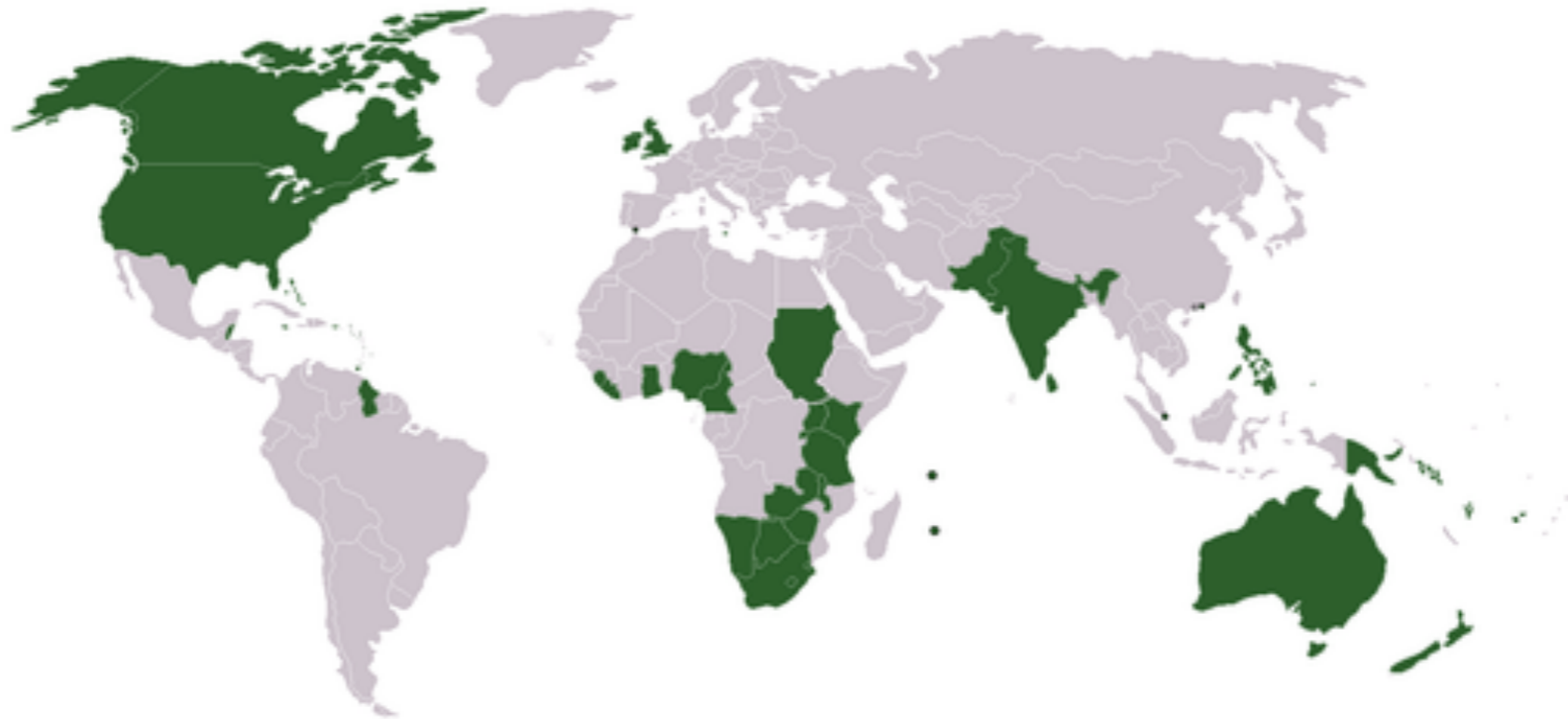
Ebenezer
@Physique_cian

Follow

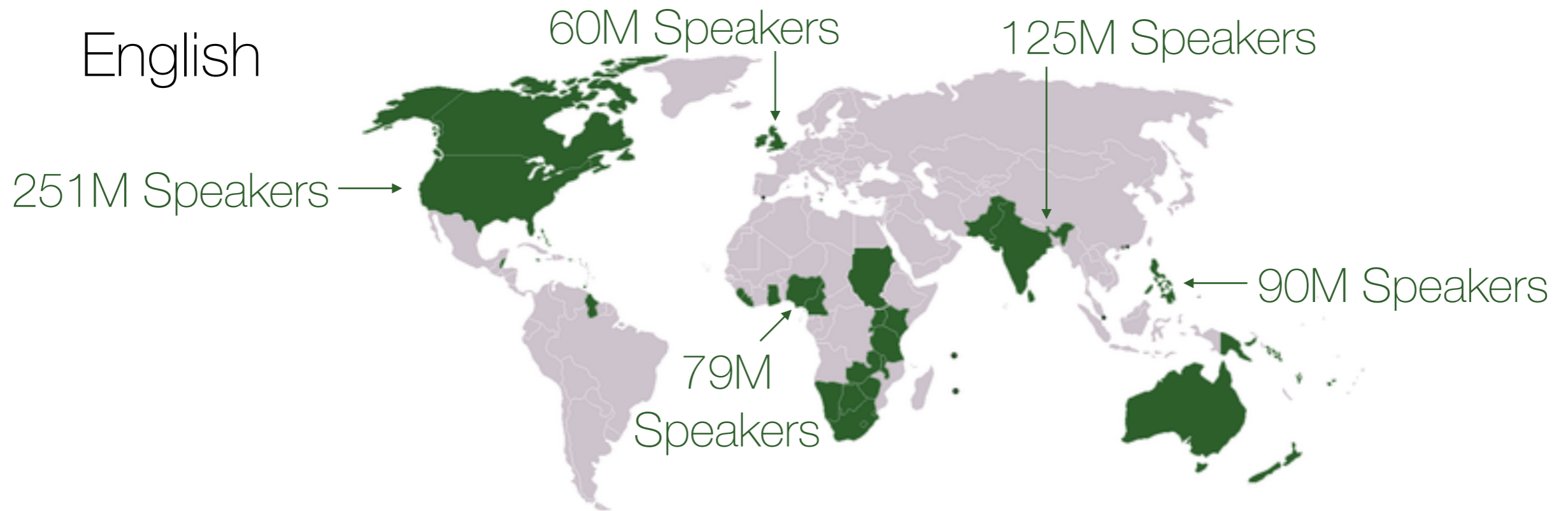
@Tblazeen R u a wizard or wat gan sef : in d mornin- u tweet, afternoon - u tweet, nyt gan u dey tweet.beta get ur IT placement wiv twitter

Global platforms attract global diversity in a language

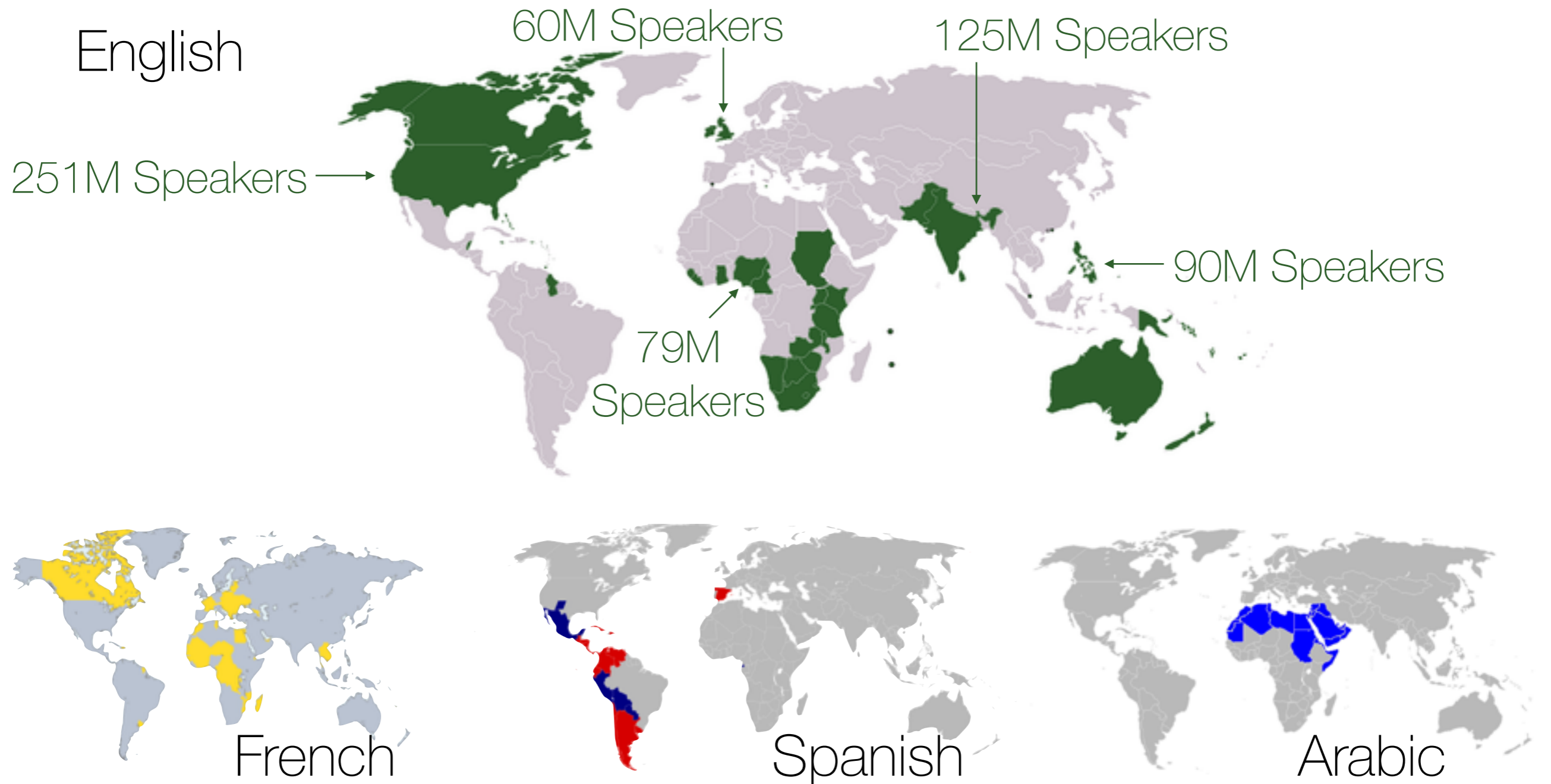
English



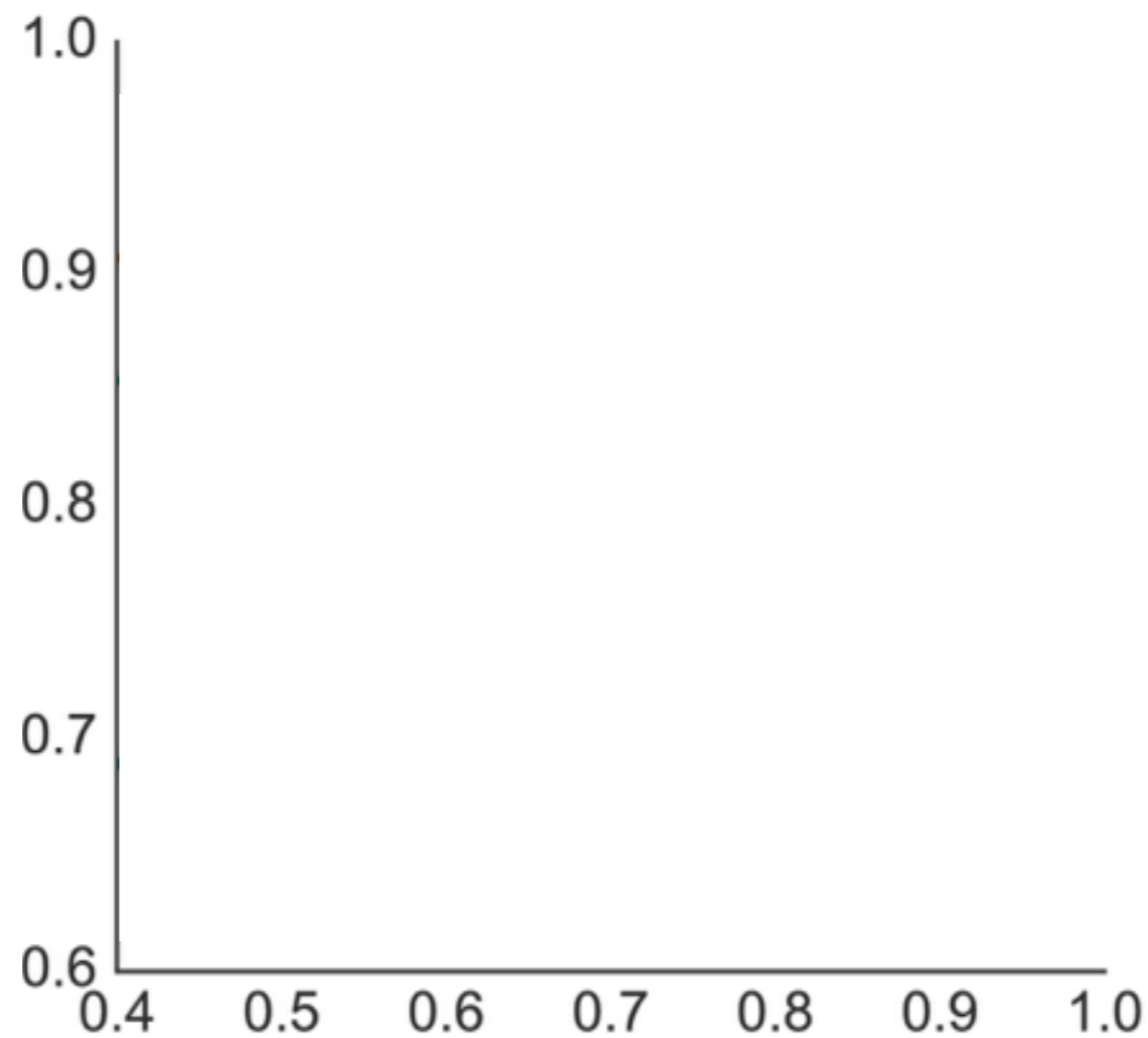
Global platforms attract global diversity in a language



Global platforms attract global diversity in a language

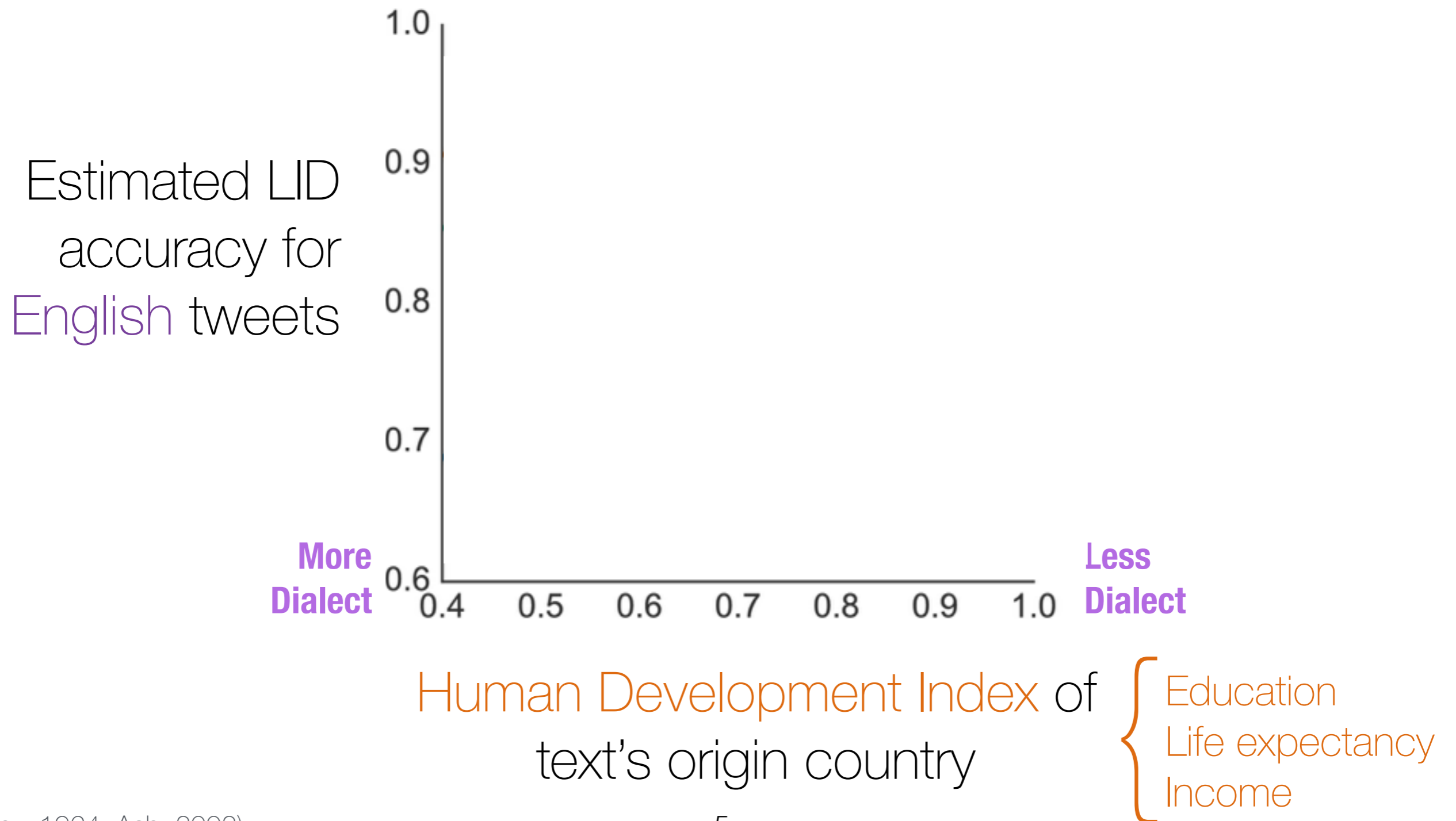


Estimated LID
accuracy for
English tweets

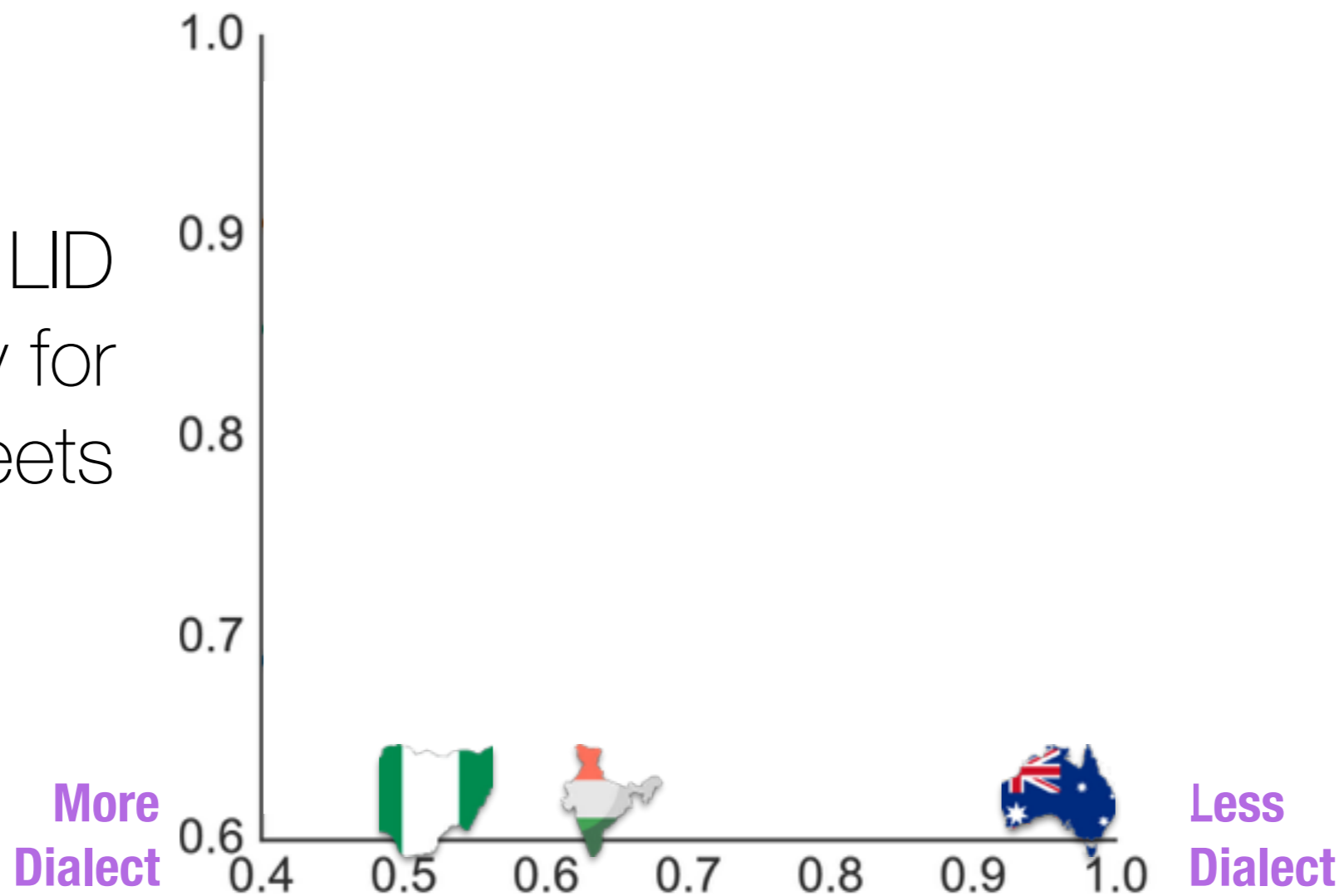


Human Development Index of
text's origin country

{ Education
Life expectancy
Income



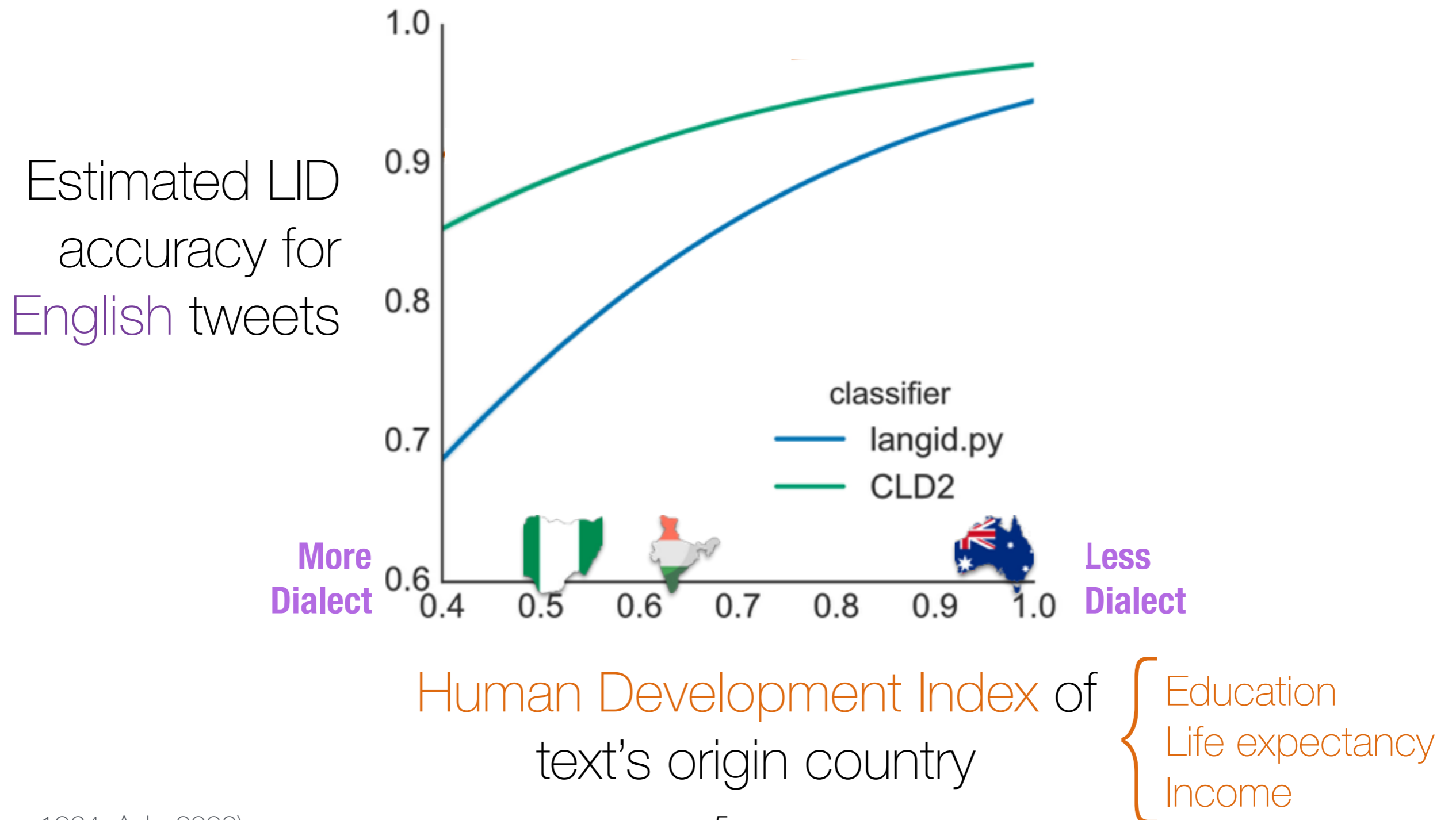
Estimated LID accuracy for English tweets



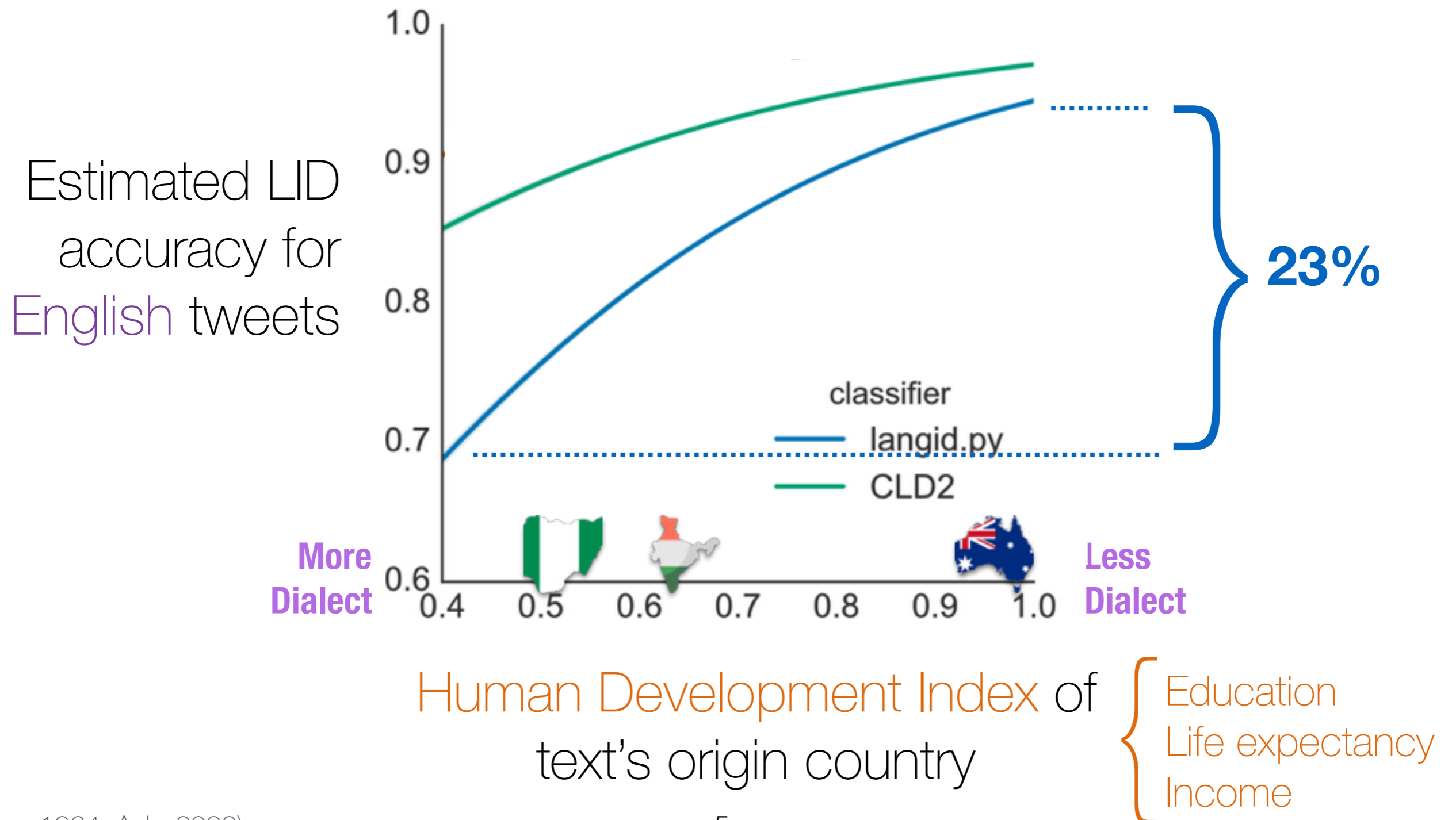
Human Development Index of text's origin country

{ Education
Life expectancy
Income

Current language detection methods perform significantly worse in less-developed countries



Current language detection methods perform significantly worse in less-developed countries



Practical Motivation: Epidemic Detection



Keyword Filter

“flu”, “sick”

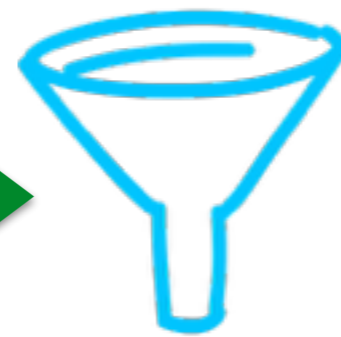
NLP

Which symptoms?

Practical Motivation: Epidemic Detection



**Language
Detection**



Keyword Filter
“flu”, “sick”

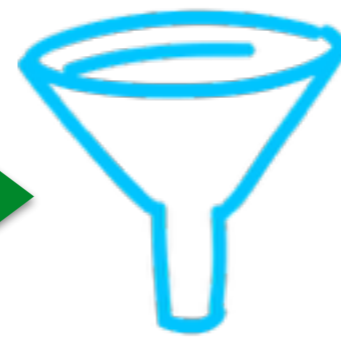


NLP
Which symptoms?

Practical Motivation: Epidemic Detection



**Language
Detection**



Keyword Filter
“flu”, “sick”



NLP
Which symptoms?

non-English



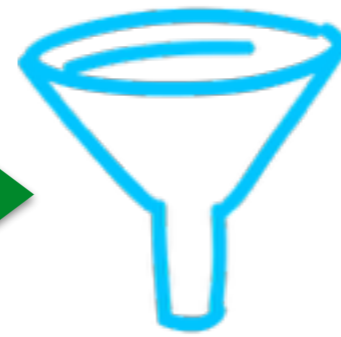
Practical Motivation: Epidemic Detection



got the flu over the weekend and I didn't know until today, & I somehow managed to give it to FIVE of my friends!!!!!!



Language Detection



Keyword Filter

"flu", "sick"



NLP

Which symptoms?



Practical Motivation: Epidemic Detection



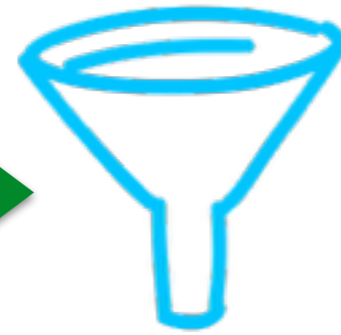
got the flu over the weekend and I didn't know until today, & I somehow managed to give it to FIVE of my friends!!!!!!



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



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



Language Detection

Keyword Filter

"flu", "sick"

NLP

Which symptoms?

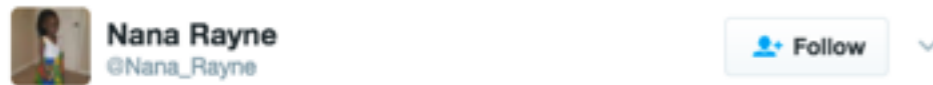
non-English



Practical Motivation: Epidemic Detection



got the flu over the weekend and I didn't know until today, & I somehow managed to give it to FIVE of my friends!!!!!!



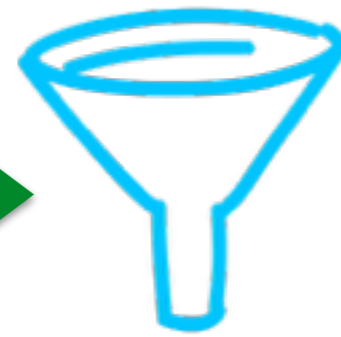
Like serious dis flu nor dey wan go oooo.... Sick



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



Language Detection



Keyword Filter

"flu", "sick"



NLP

Which symptoms?

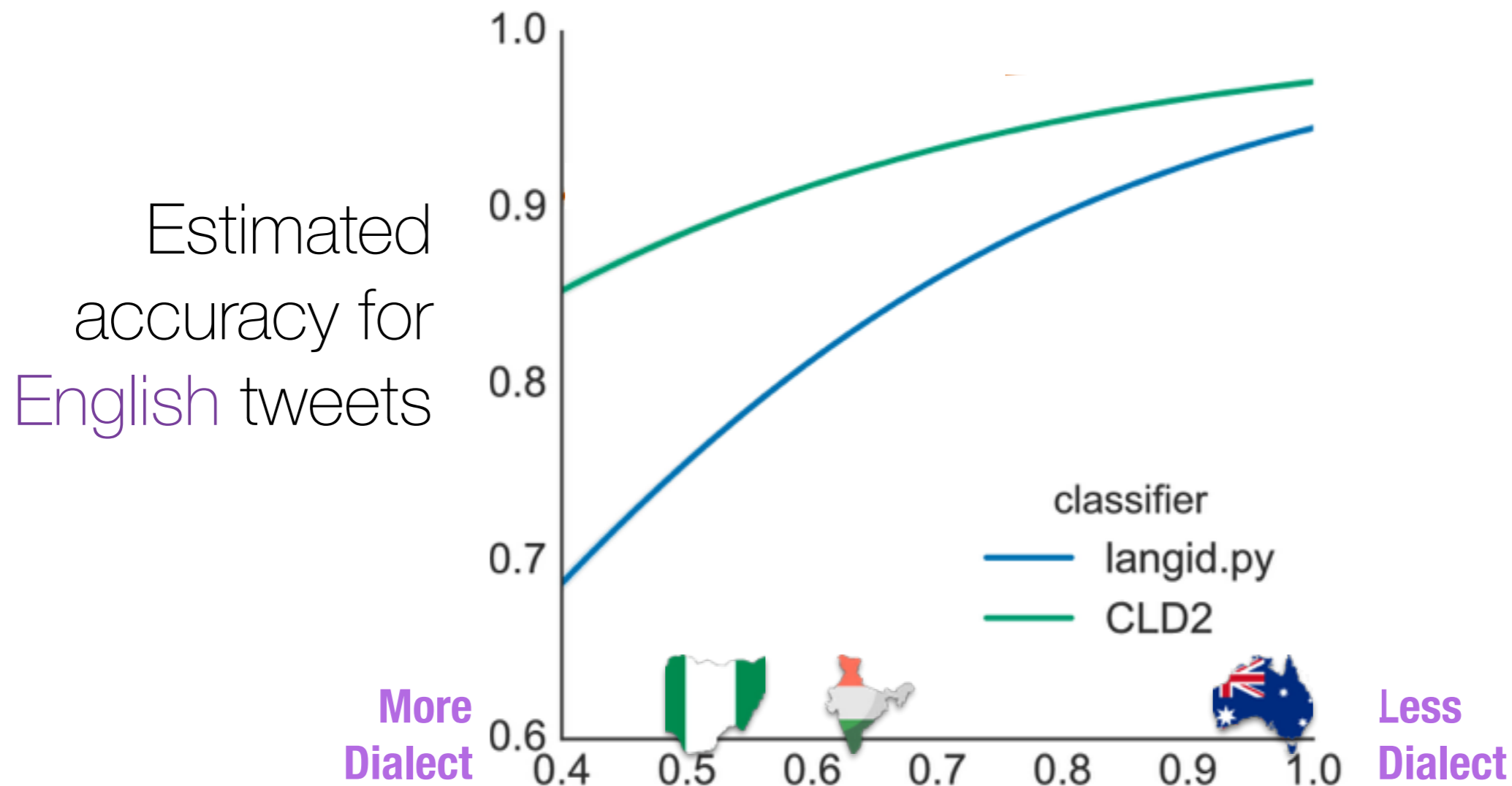
non-English?





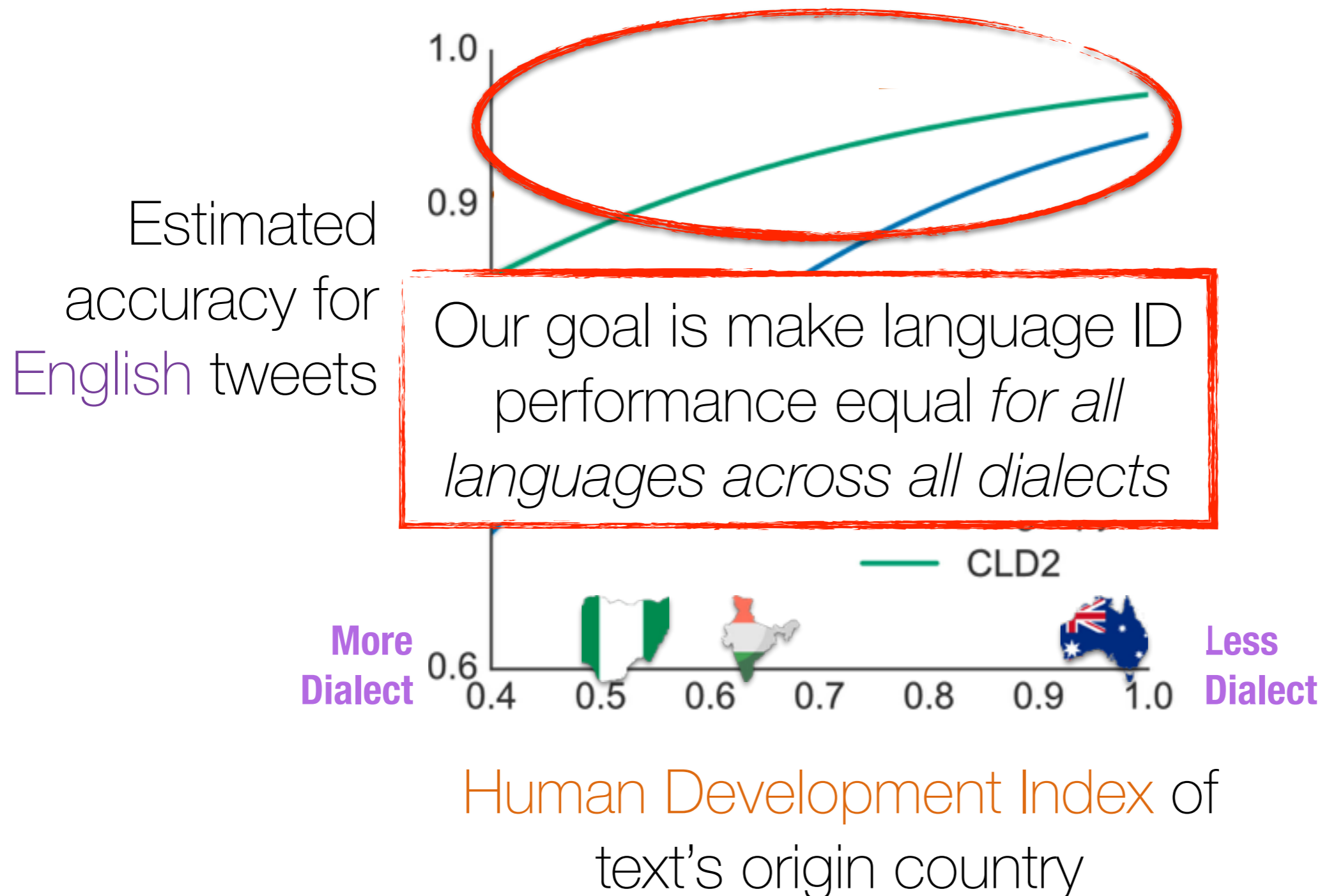
**Failing to recognize a
language silences its
speakers' voices**

Current language detection methods perform significantly worse in less-developed countries



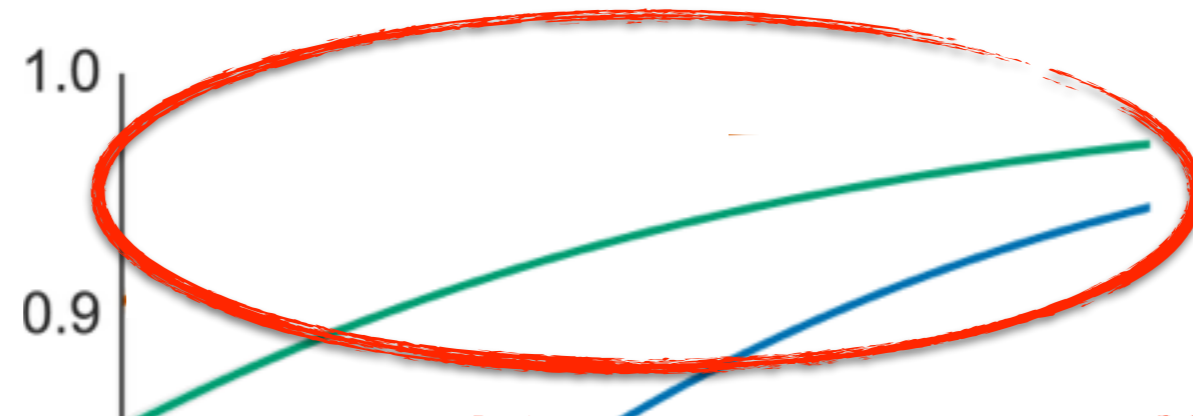
Human Development Index of text's origin country

Current language detection methods perform significantly worse in less-developed countries



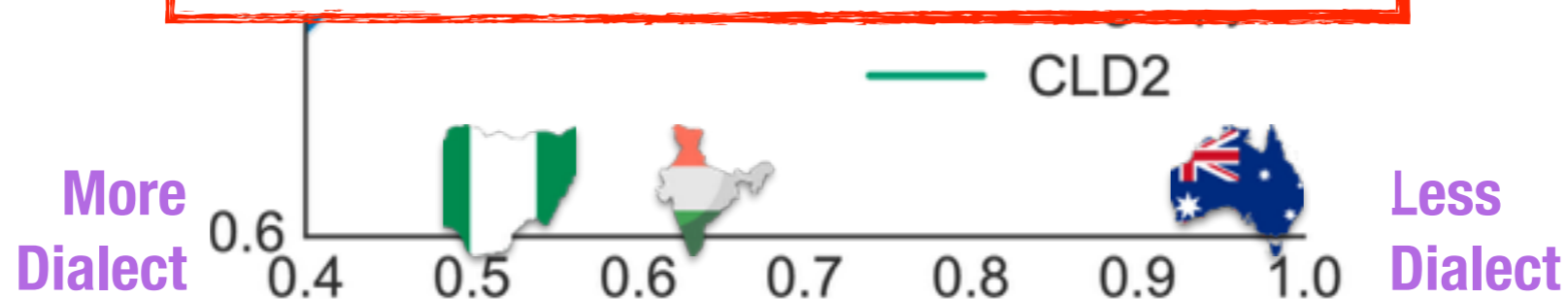
Current language detection methods perform significantly worse in less-developed countries

Estimated accuracy for English tweets



Our goal is make language ID performance equal *for all languages across all dialects*

This is a universal NLP issue!



Human Development Index of text's origin country

Key Problems: Current methods struggle in the **global** setting because

Key Problems: Current methods struggle in the **global** setting because

Data: No corpora that captures global variation in lexicon and dialect



Nana Rayne
@Nana_Rayne

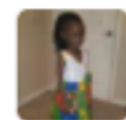
 Follow



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

Key Problems: Current methods struggle in the **global** setting because

Data: No corpora that captures global variation in lexicon and dialect



Nana Rayne
@Nana_Rayne

 Follow

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

Model: makes simplistic assumptions about how multilinguals communicate



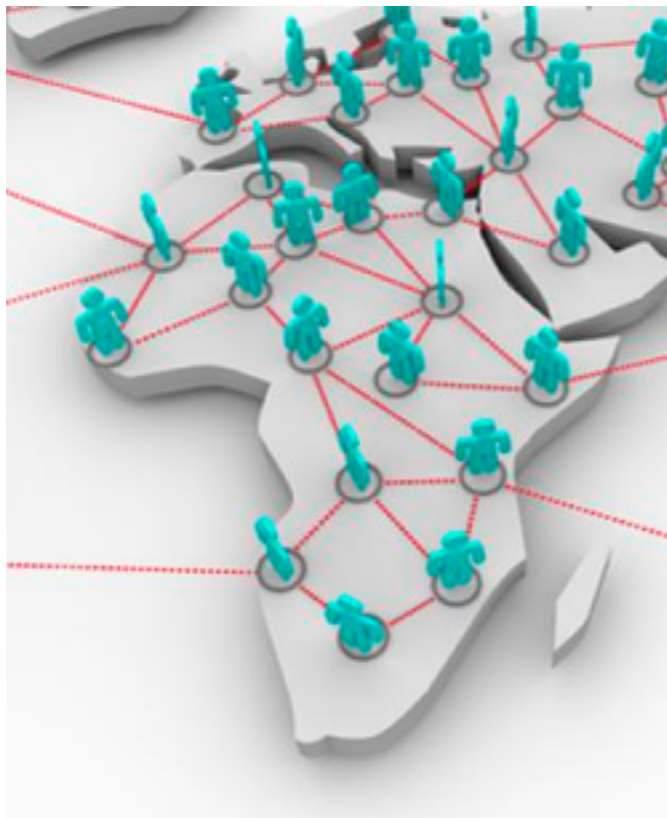
Venus
@christinedarvin

 Follow

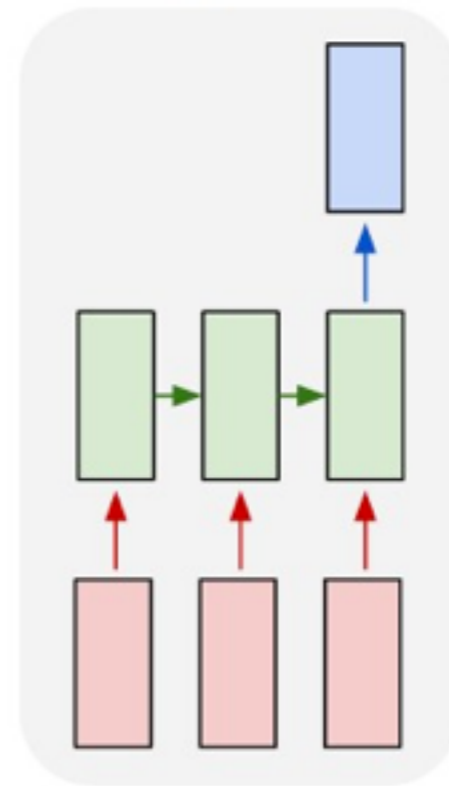
@_rkpnrnte hindi ko alam babe eh, absent ako kanina I'm sick rn hahaha 🤔👏👏

Our approach

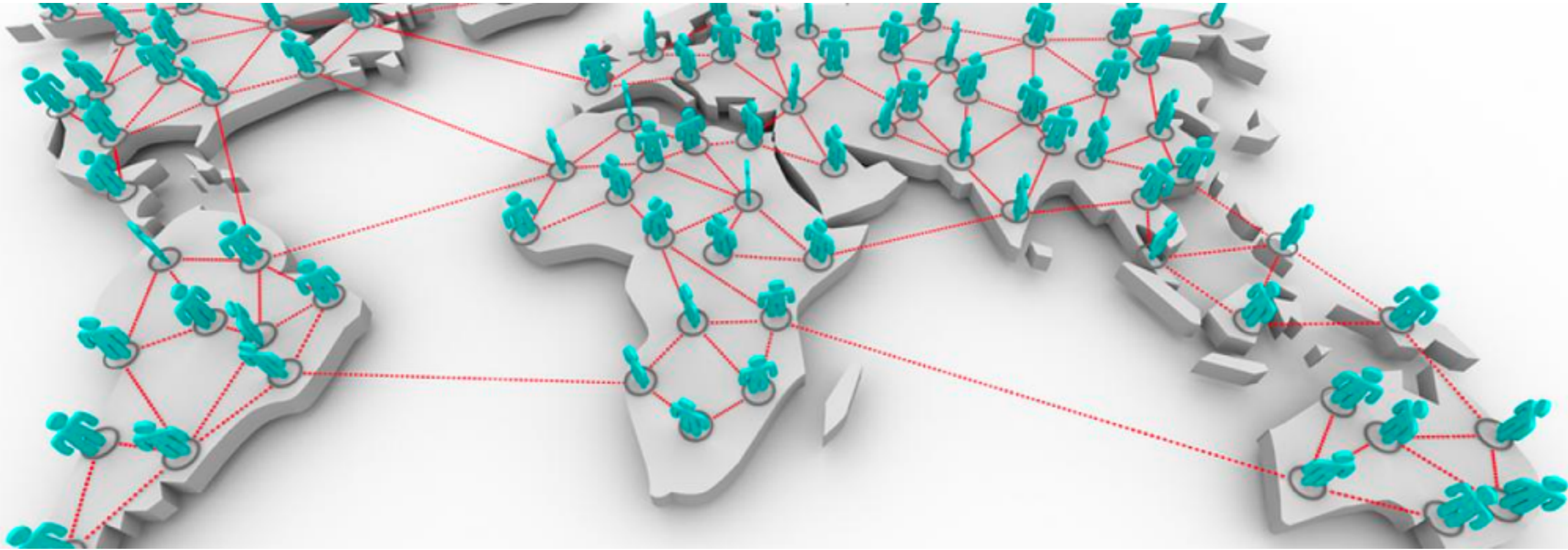
Better social representation
through network-based
sampling



NLP methodologies
capable of handling
linguistic variation



Our Data Solution: Improve linguistic representation through network-based sampling



Our Data Solution: Improve linguistic representation through network-based sampling

Bootstrap dialectic corpora using existing classifiers to find monolingual individuals



Our Data Solution: Improve linguistic representation through network-based sampling

Bootstrap dialectic corpora using existing classifiers to find monolingual individuals

*Awfully. And now,
having done wrong
to millions while
intending only good
to hundreds, I pray
God that He may
quickly take me from
a world where all I*

Our Data Solution: Improve linguistic representation through network-based sampling

Bootstrap dialectic corpora using existing classifiers to find monolingual individuals

*Amazingly. And now,
having done wrong
to millions while
intending only good
to hundreds, I pray
God that He may
quickly take me from
a world where all I*

eng

Our Data Solution: Improve linguistic representation through network-based sampling

Bootstrap dialectic corpora using existing classifiers to find monolingual individuals

*Amazingly. And now,
having done wrong
to millions while
intending only good
to hundreds, I pray
God that He may
quickly take me from
a world where all I*

eng
eng

Our Data Solution: Improve linguistic representation through network-based sampling

Bootstrap dialectic corpora using existing classifiers to find monolingual individuals

*Amazingly. And now,
having done wrong
to millions while
intending only good
to hundreds, I pray
God that He may
quickly take me from
a world where all I*

eng
eng
eng
eng
eng
eng
fra

Our Data Solution: Improve linguistic representation through network-based sampling

Bootstrap dialectic corpora using existing classifiers to find monolingual individuals

*Amazingly. And now,
having done wrong
to millions while
intending only good
to hundreds, I pray
God that He may
quickly take me from
a world where all I*

eng
eng
eng
eng
eng
eng
~~fra~~ eng

Our Data Solution: Improve linguistic representation through network-based sampling

Bootstrap dialectic corpora using existing classifiers to find monolingual individuals

*Awfully. And now,
having done wrong
to millions while
intending only good
to hundreds, I pray
God that He may
quickly take me from
a world where all I*

eng
eng
eng
eng
eng
eng
~~fra~~ eng

Sample from the geolocated Twitter social network to include text from people at all locations

Build a strategically-diverse corpora

Build a strategically-diverse corpora

Topical



Build a strategically-diverse corpora

Topical



Geographic



Build a strategically-diverse corpora

Topical



Geographic



Social



Build a strategically-diverse corpora and synthesize code-switched examples

Topical



Geographic



Social



Multilingual



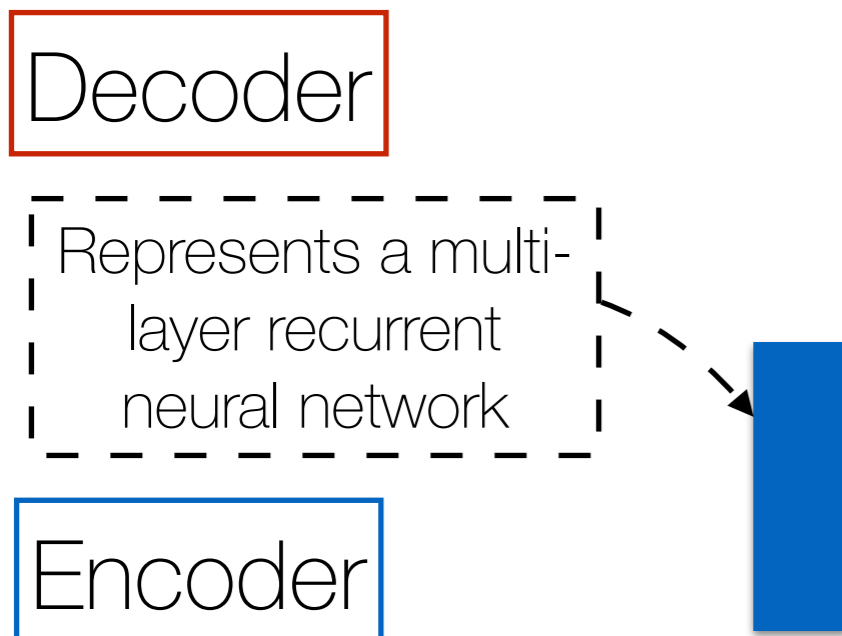
Our model solution: treat language identification as a **character-based sequence to sequence** task.

Decoder

Encoder

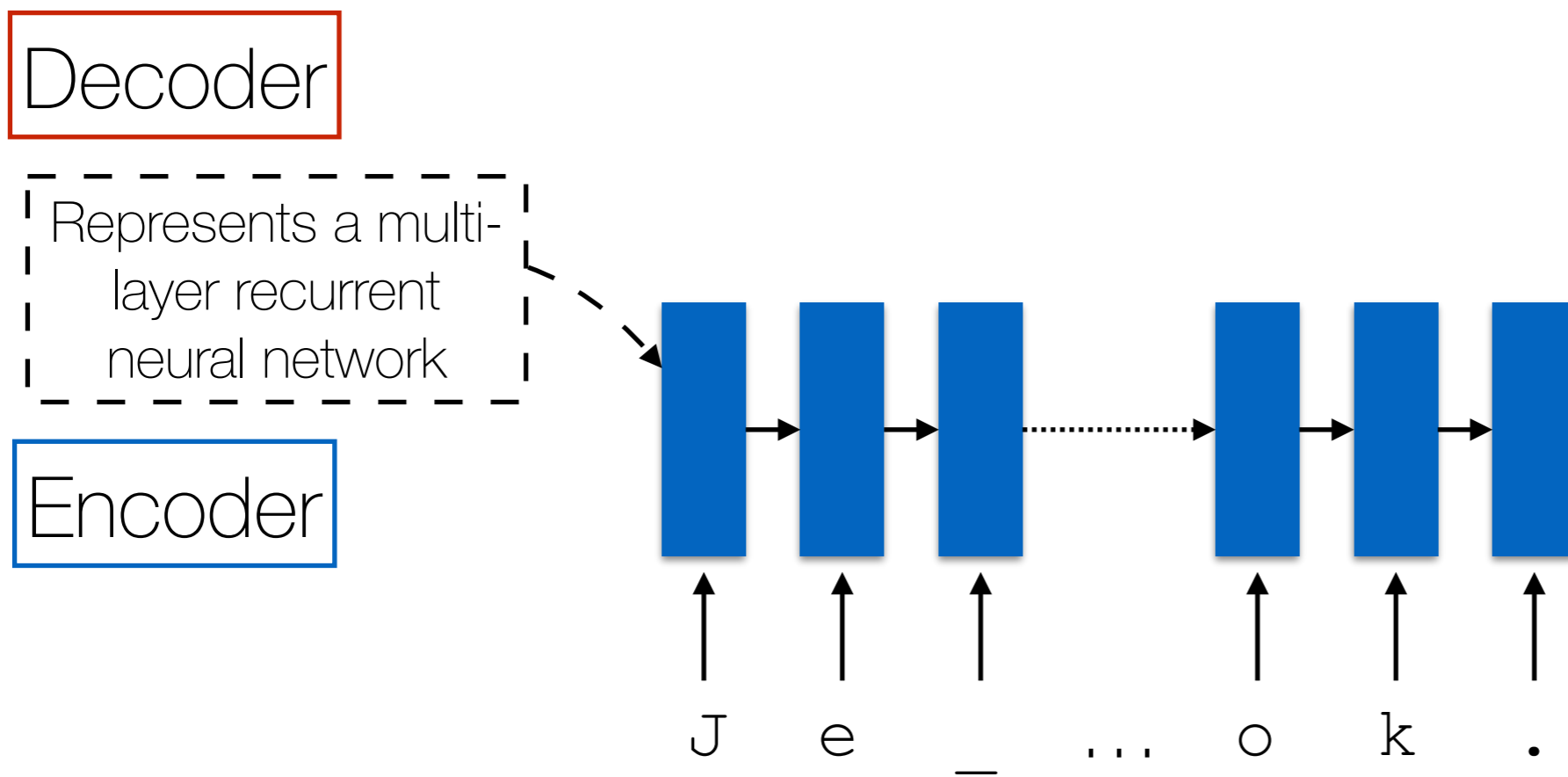
Je vais commander à emporter. I'm too lazy to cook.

Our model solution: treat language identification as a **character-based sequence to sequence** task.



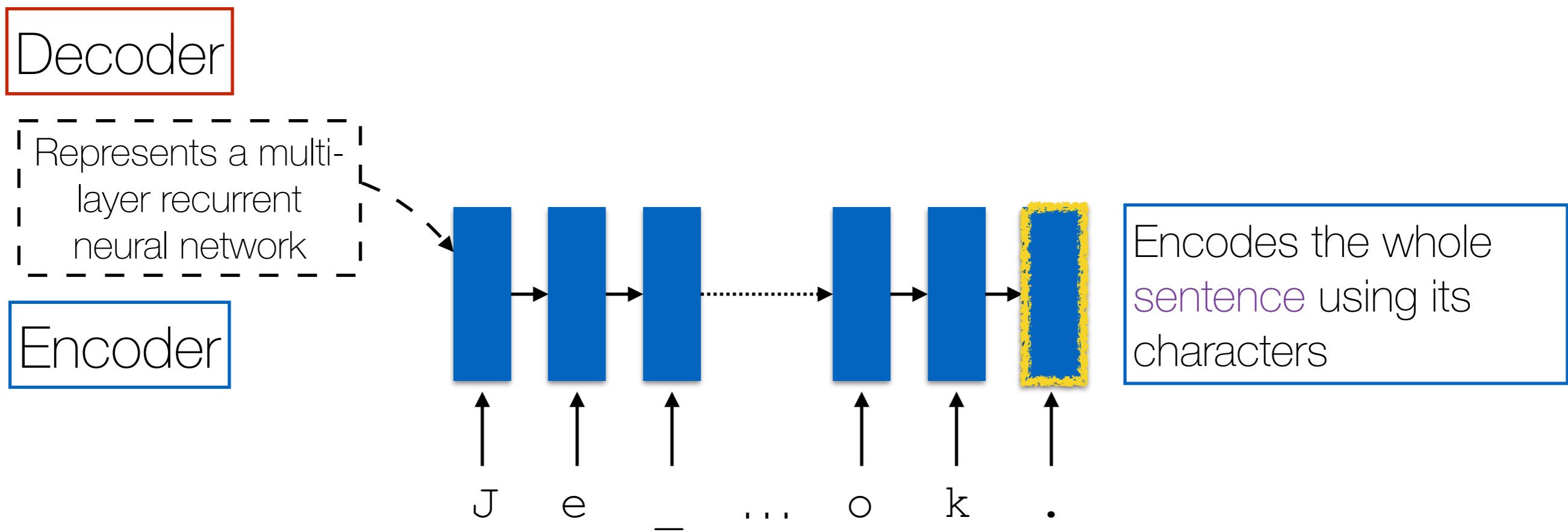
Je vais commander à emporter. I'm too lazy to cook.

Our model solution: treat language identification as a **character-based sequence to sequence** task.



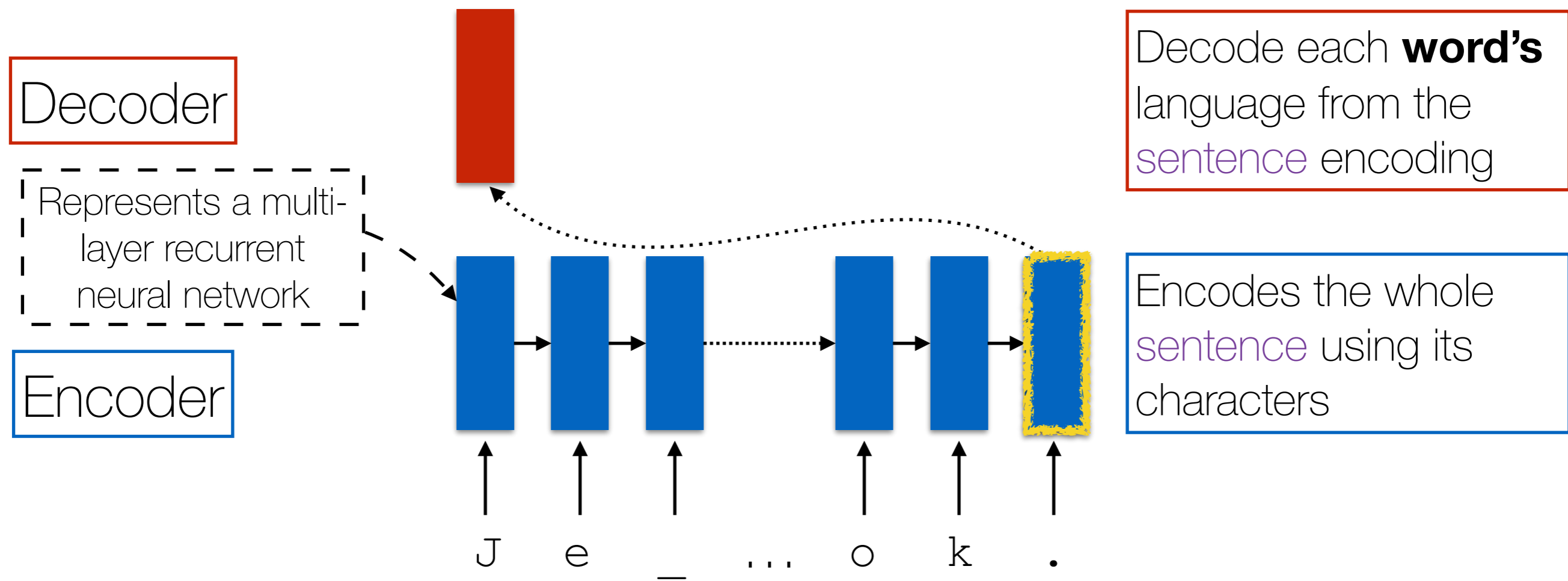
Je vais commander à emporter. I'm too lazy to cook.

Our model solution: treat language identification as a **character-based sequence to sequence** task.



Je vais commander à emporter. I'm too lazy to cook.

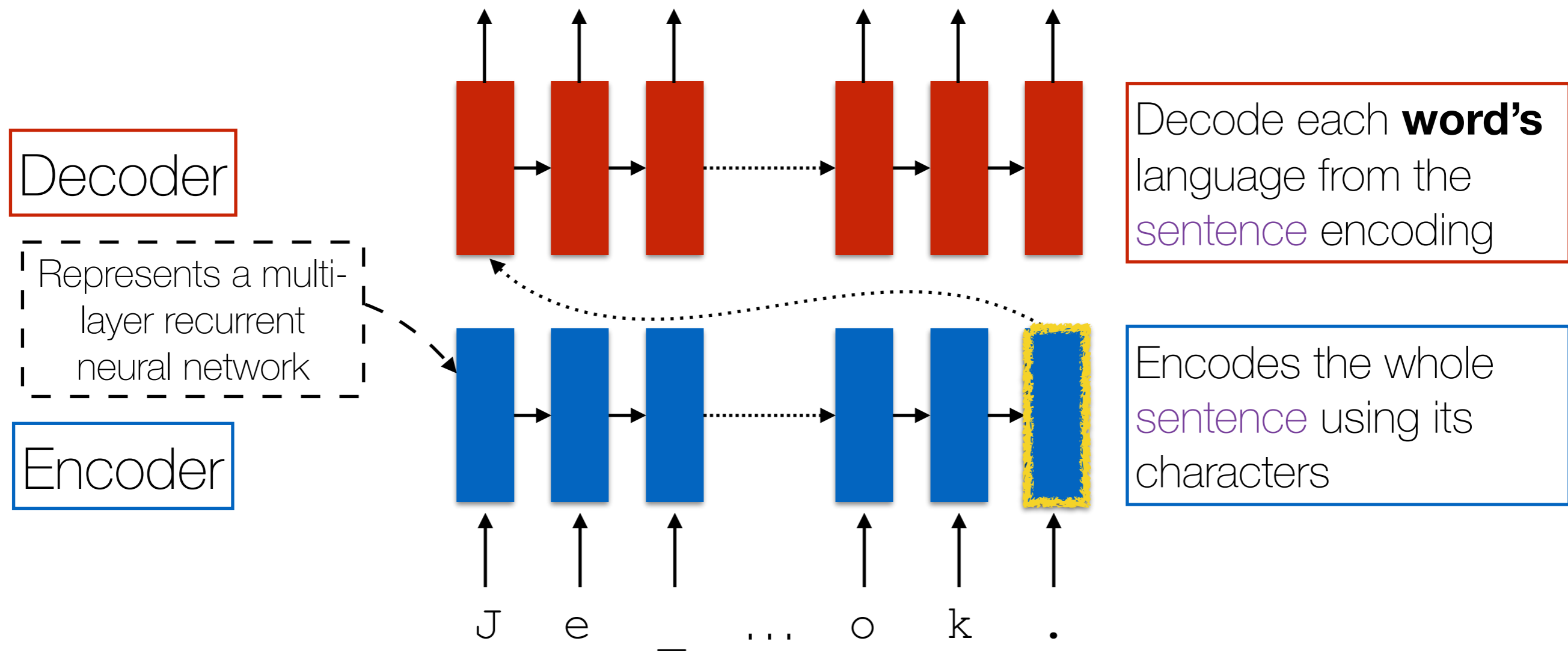
Our model solution: treat language identification as a **character-based sequence to sequence** task.



Je vais commander à emporter. I'm too lazy to cook.

Our model solution: treat language identification as a **character-based sequence to sequence** task.

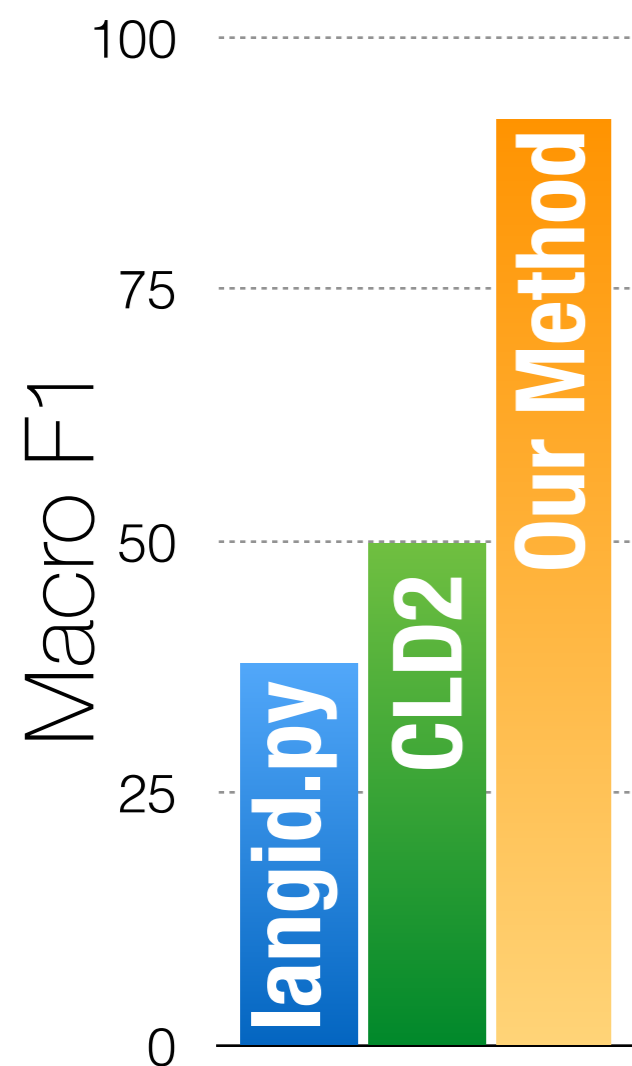
Fra Fra Fra Fra Fra . Eng Eng Eng Eng Eng .



Je vais commander à emporter. I'm too lazy to cook.

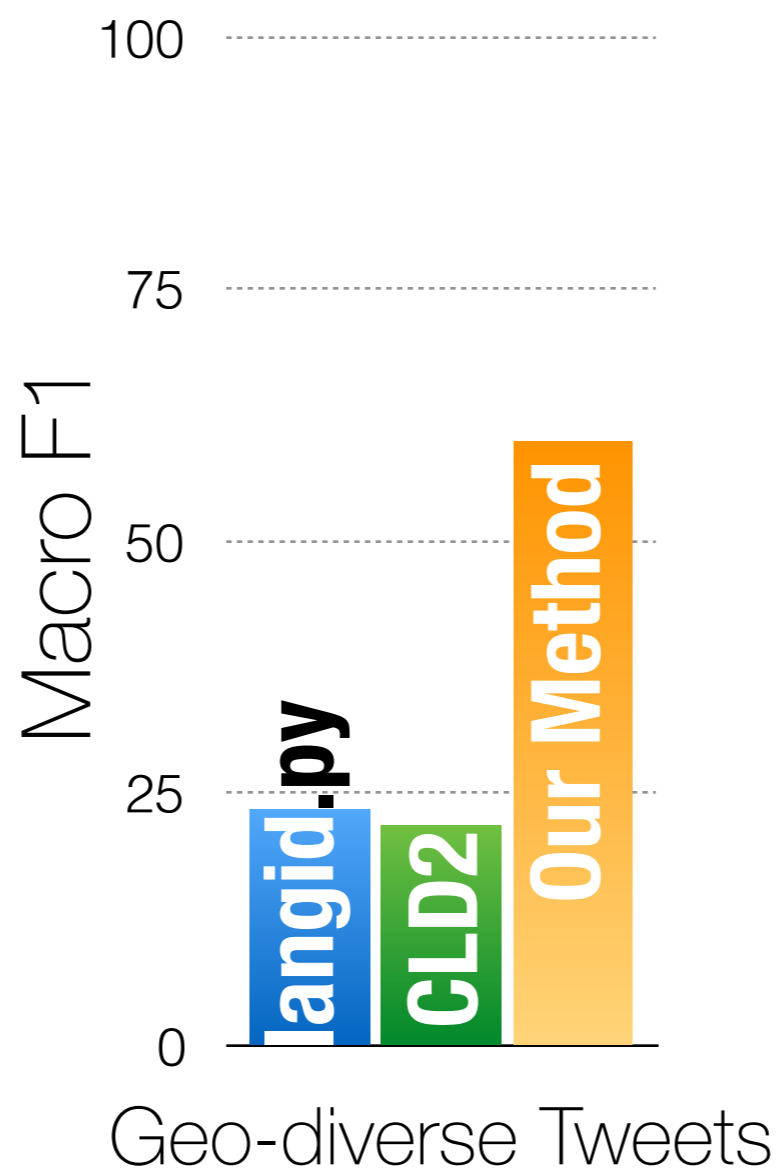
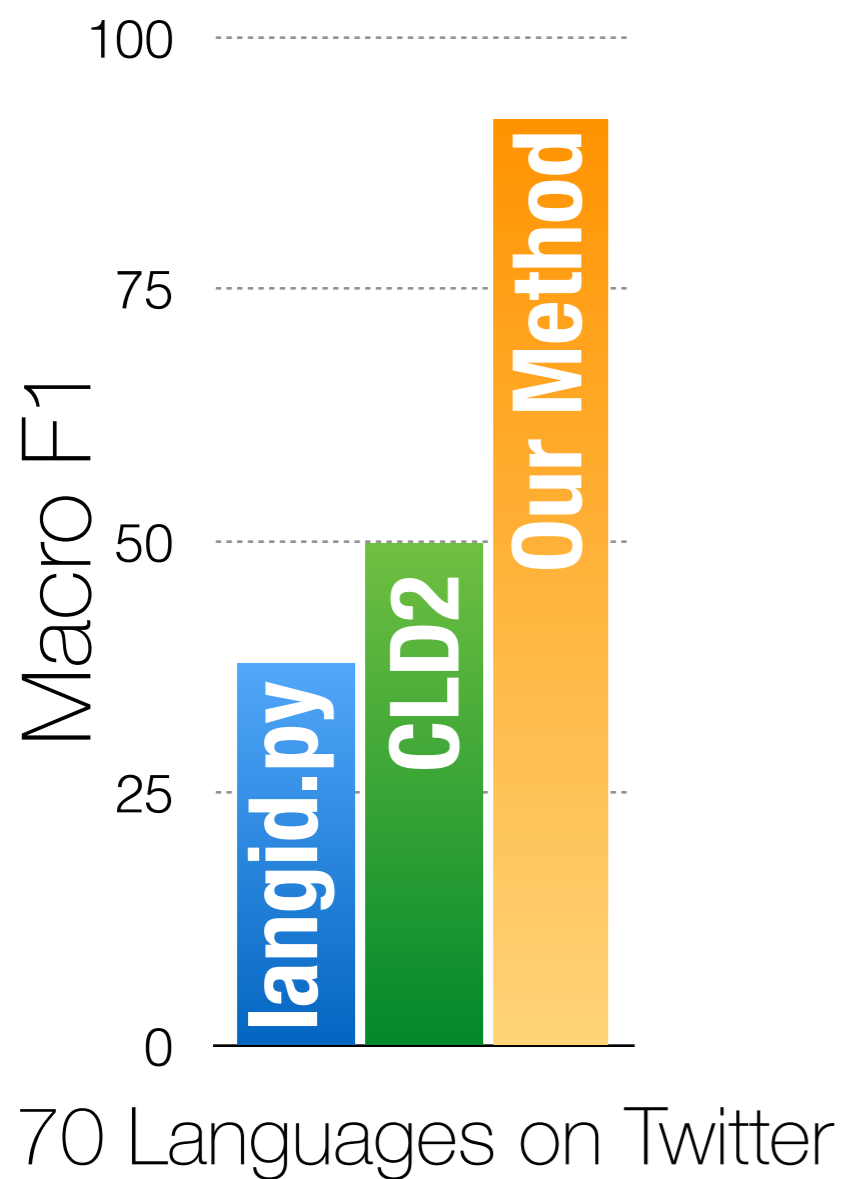
Equilid vs off-the-shelf

Equilid vs off-the-shelf

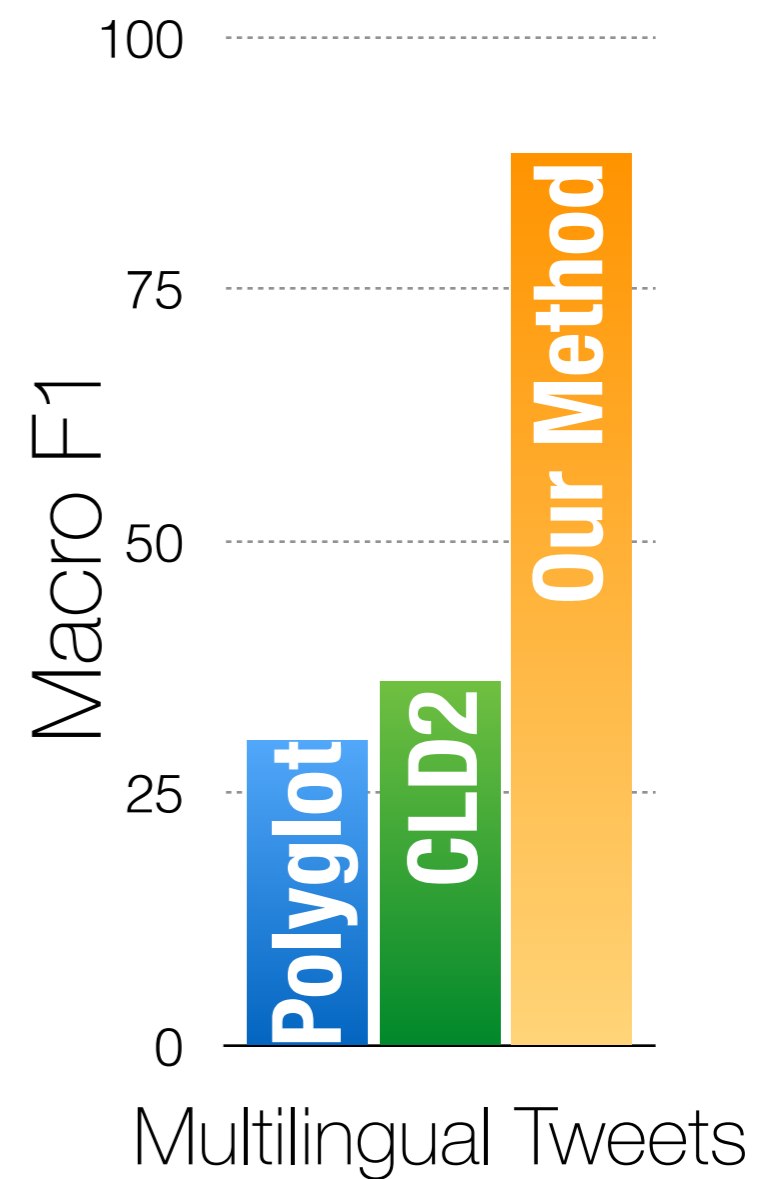
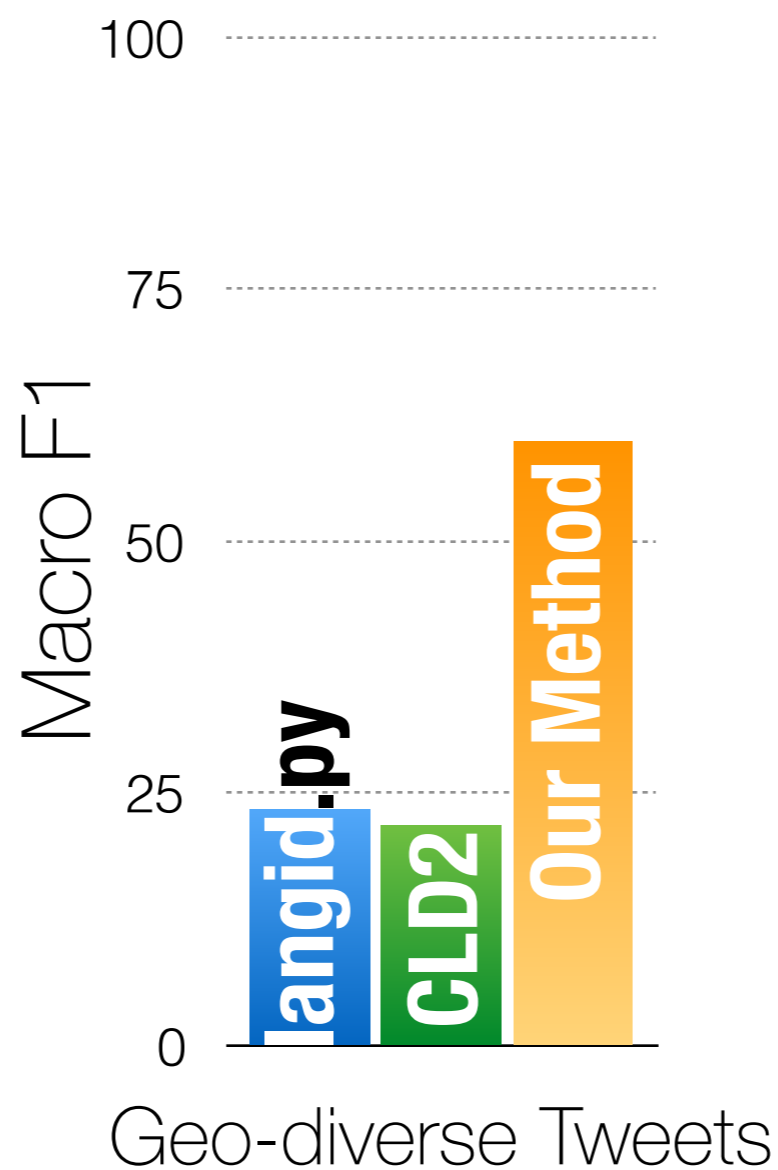
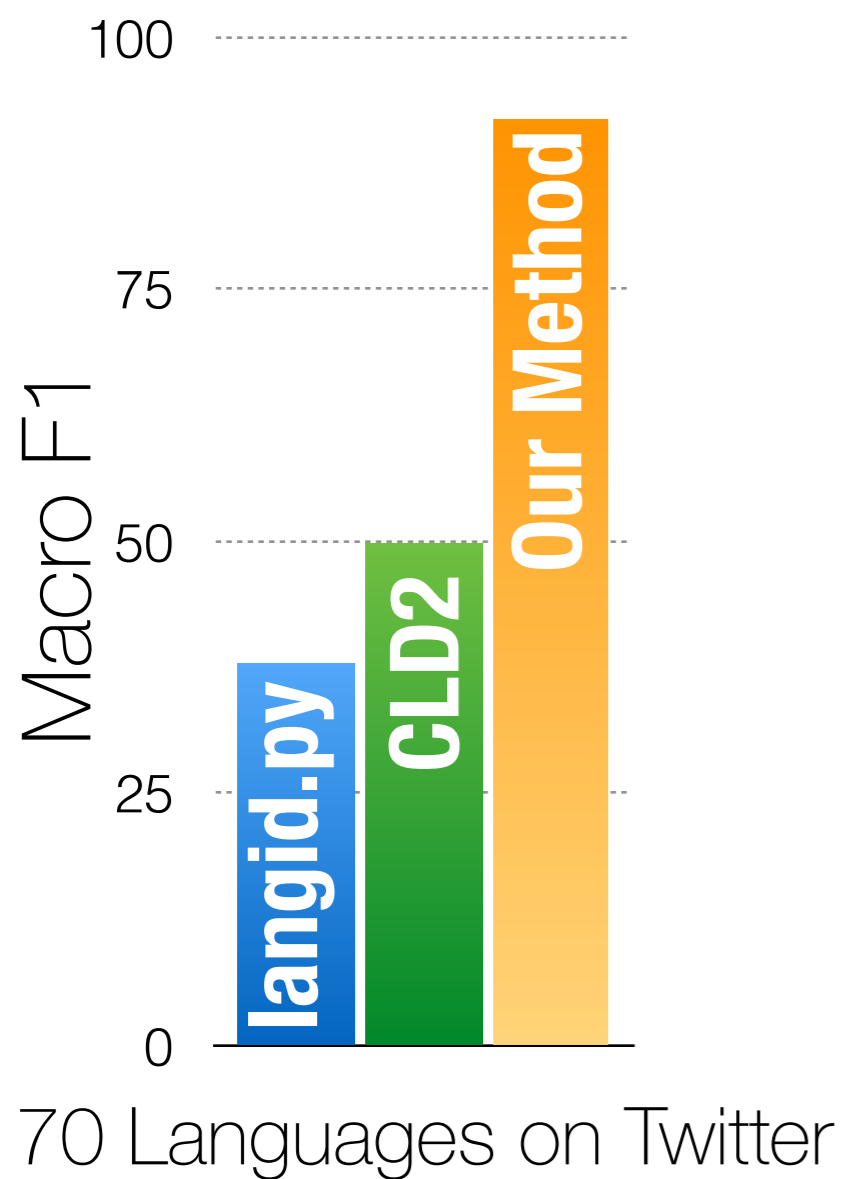


70 Languages on Twitter

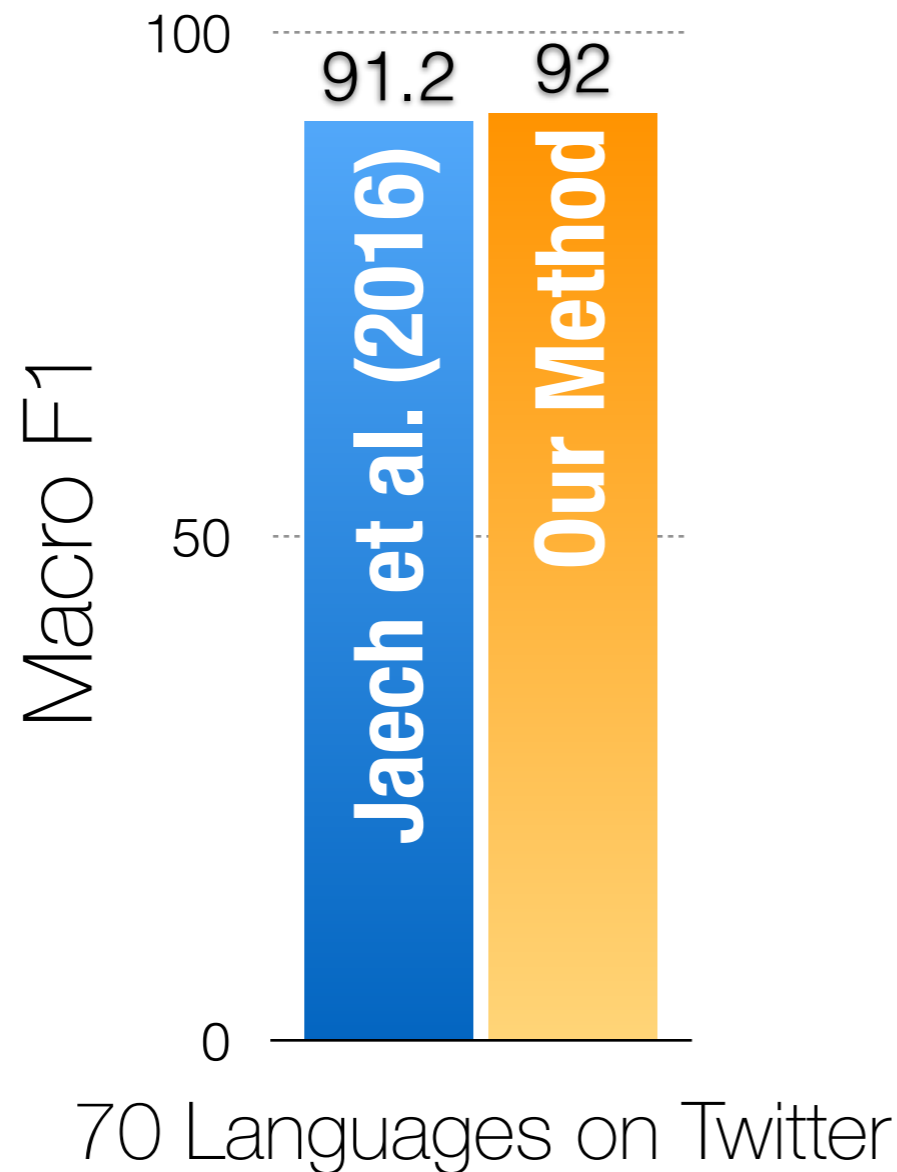
Equilid vs off-the-shelf



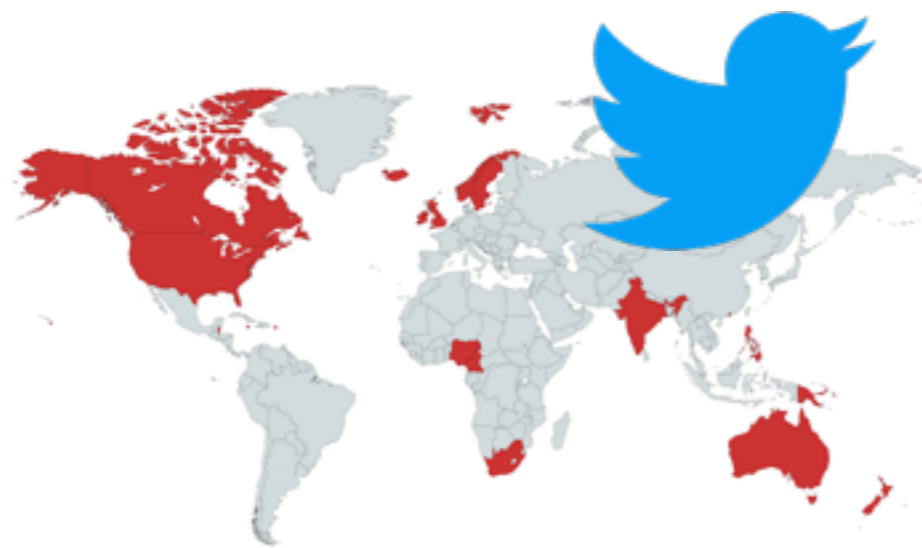
Equilid vs off-the-shelf



Equilid even outperforms system specifically tuned for each dataset

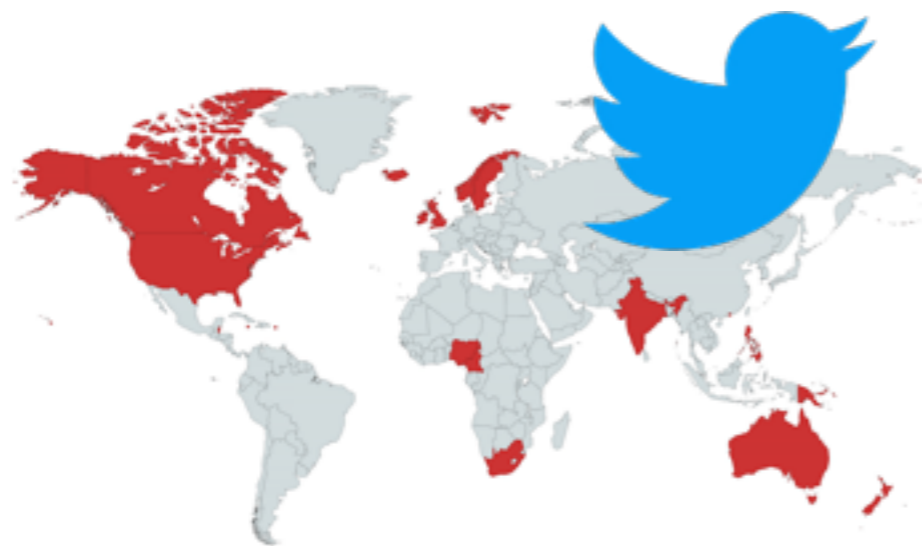


Case Study: Do our solutions provide socially-equitable language identification for health-related queries?



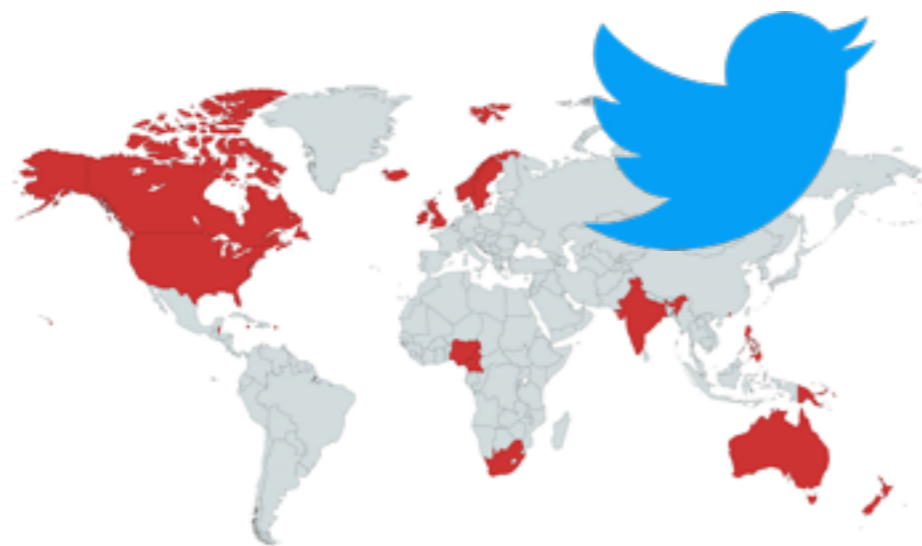
1M Tweets with any of 385
English terms from
established lexicons for
influenza, psychological well-
being, and social health

Case Study: Do our solutions provide socially-equitable language identification for health-related queries?



1M Tweets with any of 385
English terms from
established lexicons for
influenza, psychological well-
being, and social health

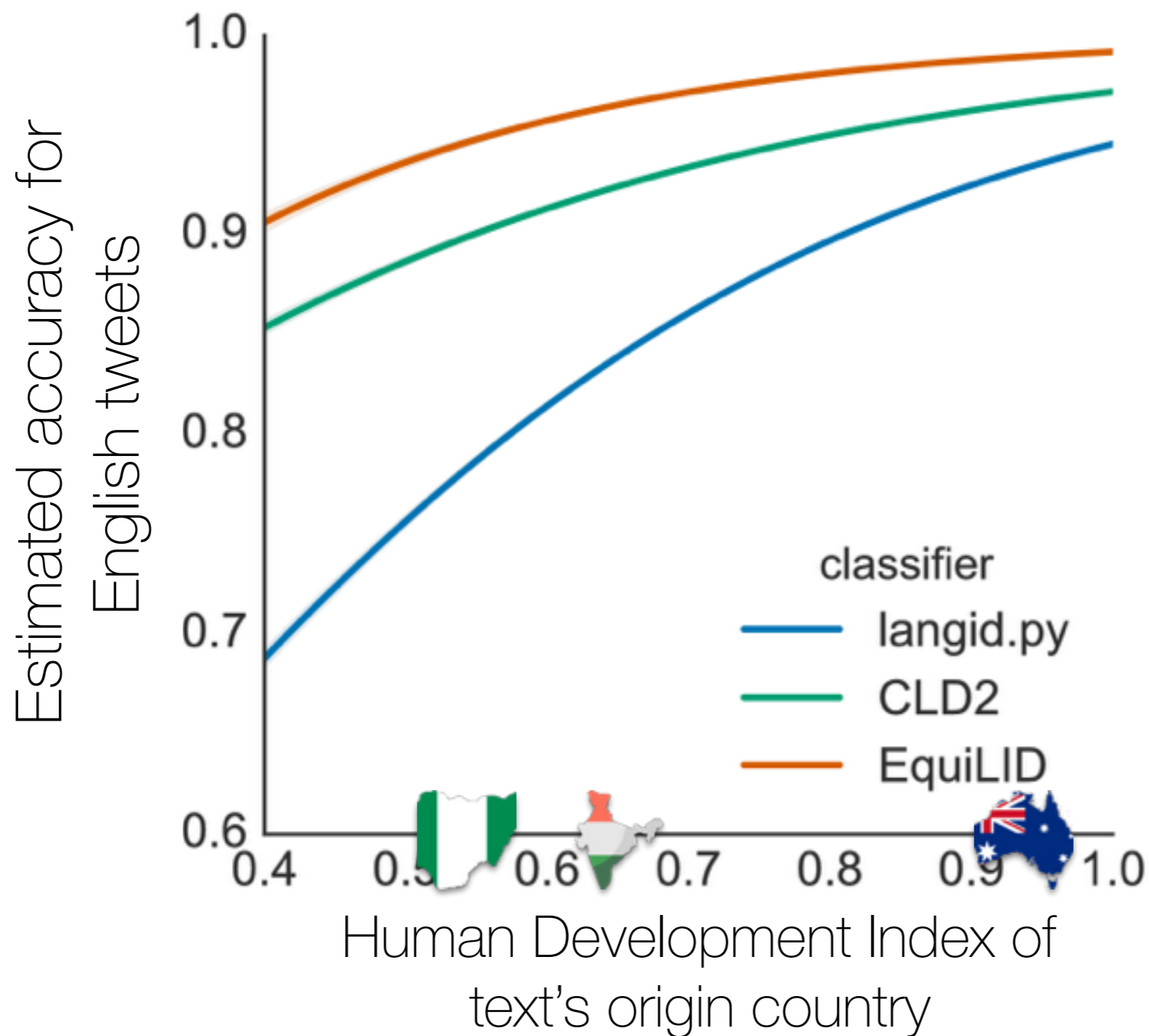
Case Study: Do our solutions provide socially-equitable language identification for health-related queries?



1M Tweets with any of 385
English terms from
established lexicons for
influenza, psychological well-
being, and social health

Task: does the language identification system recognize every tweet as English?

Equilid raises the bar for socially-equitable language identification



Social Equality doesn't stop at Language Identification



Better social
representation in
our data

Methodologies
capable of handling
language as it is used

Social Equality doesn't stop at Language Identification



Better social
representation in
our data



Methodologies
capable of handling
language as it is used

Be equitable!

<https://github.com/davidjurgens/equilid>

David Jurgens, Yulia Tsvetkov, and Dan Jurafsky

