

# Geographically-Informed Language Identification

Jonathan Dunn

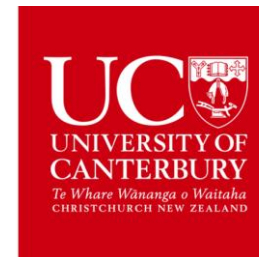
[jedunn@illinois.edu](mailto:jedunn@illinois.edu)

[www.jdunn.name](http://www.jdunn.name)



Lane Edwards-Brown

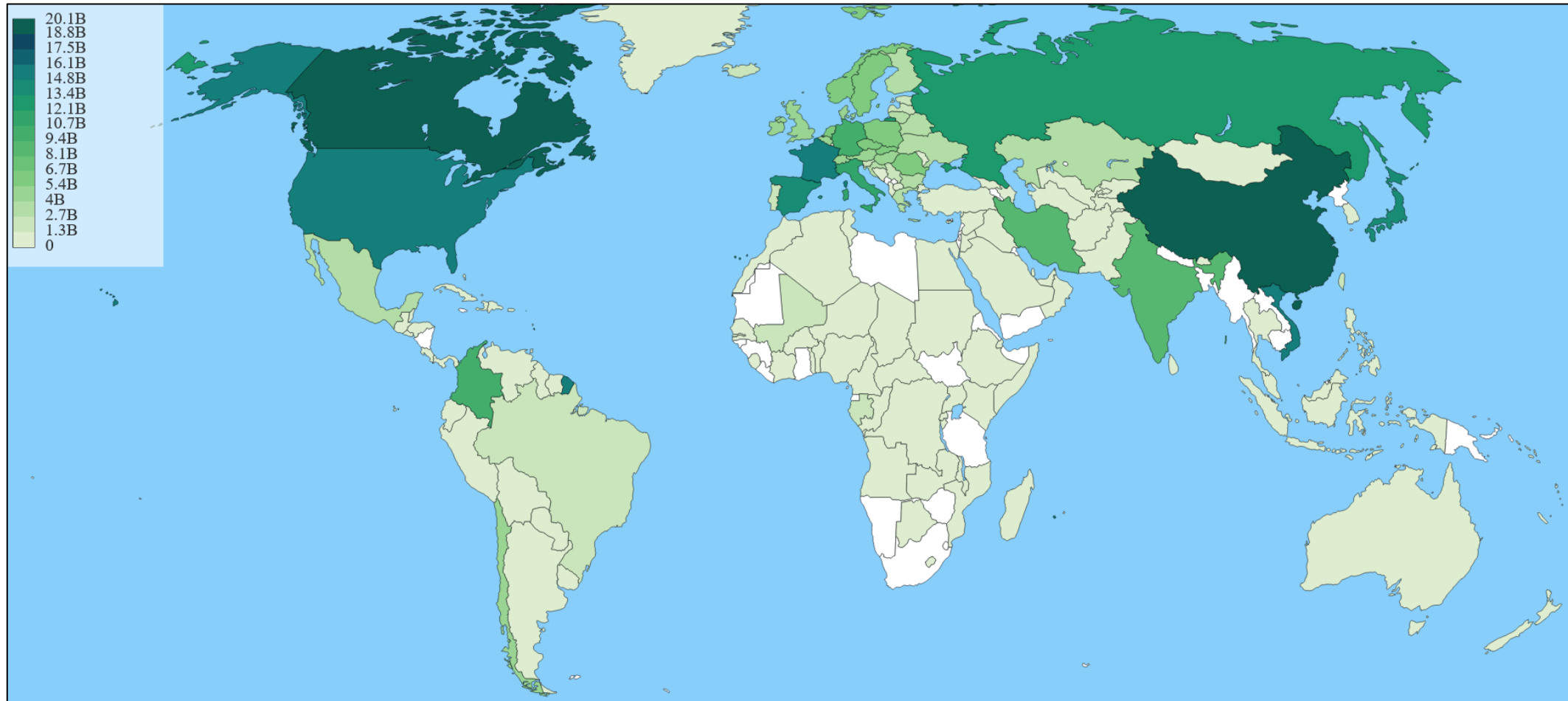
[laneedwardsbrown@gmail.com](mailto:laneedwardsbrown@gmail.com)



Do geographic priors improve LID for lower-resource languages?

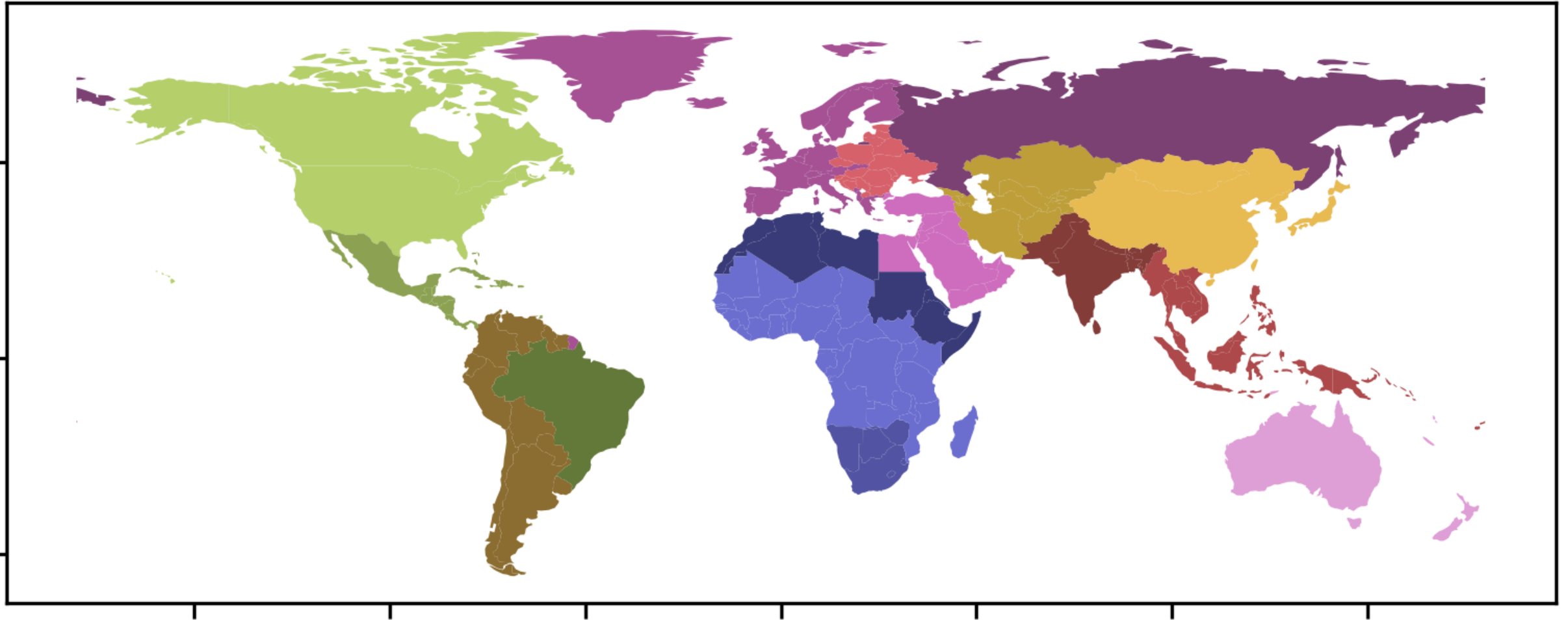
Goal: Enhanced LID for geo-referenced corpora

[www.earthLings.io](http://www.earthLings.io)



## 16 Region-Specific LID Models

---



# Data from LID Sources

---

<b>Corpus</b>	<b>N. Langs</b>
Bible Translations (Brown, 2014)	614
Global Voices News (Tiedemann, 2012)	41
JW 300 (Agić and Vulić, 2019)	380
Open Subtitles (Lison and Tiedemann, 2016)	62
QCRI Educational Domain (Tiedemann, 2012)	42
Tatoeba Sentences (Tiedemann, 2012)	309
Wikipedia Articles TensorFlow DataSets	280

Table 1: Primary Sources of Training Data

# Language Data from Glottolog

---



# International Languages

---

<b>Language</b>	<b>Abbrev.</b>	<b>Language</b>	<b>Abbrev.</b>
Amharic	amh	Korean	kor
Arabic	ara	Mandarin	zho
Bengali	ben	Marathi	mar
English	eng	Polish	pol
Farsi	fas	Portuguese	por
French	fra	Punjabi	pan
German	deu	Russian	rus
Gujarati	guj	Spanish	spa
Hausa	hau	Swahili	swa
Hindi	hin	Tagalog	tgl
Indonesian	ind	Tamil	tam
Italian	ita	Telugu	tel
Japanese	jpn	Thai	tha
Javanese	jav	Turkish	tur
Kannada	kan	Urdu	urd
		Vietnamese	vie

Table 2: International languages which are included in each regional model.

# Results

---

Region	N. Langs	F-Score		Test Samples
		<i>Geo</i>	<i>Baseline</i>	
Africa, North	44	<b>0.990</b>	0.886	621k
Africa, Southern	58	<b>0.982</b>	0.903	1,053k
Africa, Sub	166	<b>0.980</b>	0.947	1,931k
America, Central	188	<b>0.991</b>	0.961	2,965k
America, North	68	<b>0.993</b>	0.902	1,017k
America, South	129	<b>0.995</b>	0.960	4,612k
America, Brazil	88	<b>0.996</b>	0.945	818k
Asia, Central	54	<b>0.988</b>	0.906	777k
Asia, East	46	<b>0.990</b>	0.892	1,131k
Asia, South	60	<b>0.986</b>	0.914	979k
Asia, Southeast	325	<b>0.990</b>	0.972	3,992k
Europe, East	65	<b>0.978</b>	0.908	3,132k
Europe, West	108	<b>0.967</b>	0.921	5,473k
Europe, Russia	65	<b>0.984</b>	0.914	1,098k
Middle East	53	<b>0.988</b>	0.904	801k
Oceania	49	<b>0.984</b>	0.890	745k

Geo = GeoLID in current region

Baseline = Same architecture, no geo

Models = fastText-based classifier



## Results for only local languages

---

Region	N. Langs	F-Score		Test Samples
		<i>Geo</i>	<i>Baseline</i>	
Africa, North	13	<b>0.986</b>	0.972	75k
Africa, Southern	27	<b>0.976</b>	0.964	442k
Africa, Sub	135	<b>0.979</b>	0.970	1,313k
America, Central	157	<b>0.991</b>	0.983	2,339k
American, North	37	<b>0.994</b>	0.945	363k
America, South	98	<b>0.997</b>	0.994	3,747k
America, Brazil	57	<b>0.999</b>	0.996	287k
Asia, Central	23	0.984	0.983	237k
Asia, East	15	0.993	0.980	410k
Asia, South	29	0.983	0.983	340k
Asia, Southeast	294	<b>0.991</b>	0.985	3,291k
Europe, East	34	<b>0.975</b>	0.961	2,363k
Europe, West	77	<b>0.963</b>	0.950	4,565k
Europe, Russia	34	<b>0.979</b>	0.972	525k
Middle East	22	<b>0.986</b>	0.981	255k
Oceania	18	<b>0.978</b>	0.959	168k

Geo = GeoLID in current region

Baseline = Same architecture, no geo

Models = fastText-based classifier

# OpenLID as extra test set

---

Region	N.	Geo	Baseline
Africa, North	35	<b>0.990</b>	0.969
Africa, Southern	44	<b>0.988</b>	0.968
Africa, Sub	65	<b>0.983</b>	0.968
America, Central	36	<b>0.987</b>	0.970
America, North	31	<b>0.986</b>	0.966
America, South	35	<b>0.990</b>	0.969
America, Brazil	33	<b>0.993</b>	0.967
Asia, Central	43	<b>0.991</b>	0.974
Asia, East	38	0.964	0.945
Asia, South	42	0.947	0.925
Asia, Southeast	48	0.936	0.926
Europe, East	53	<b>0.992</b>	0.978
Europe, West	67	<b>0.991</b>	0.982
Europe, Russia	42	<b>0.992</b>	0.974
Middle East	41	<b>0.994</b>	0.973
Oceania	35	0.976	0.960

Geo = GeoLID in current region

Baseline = Same architecture, no geo

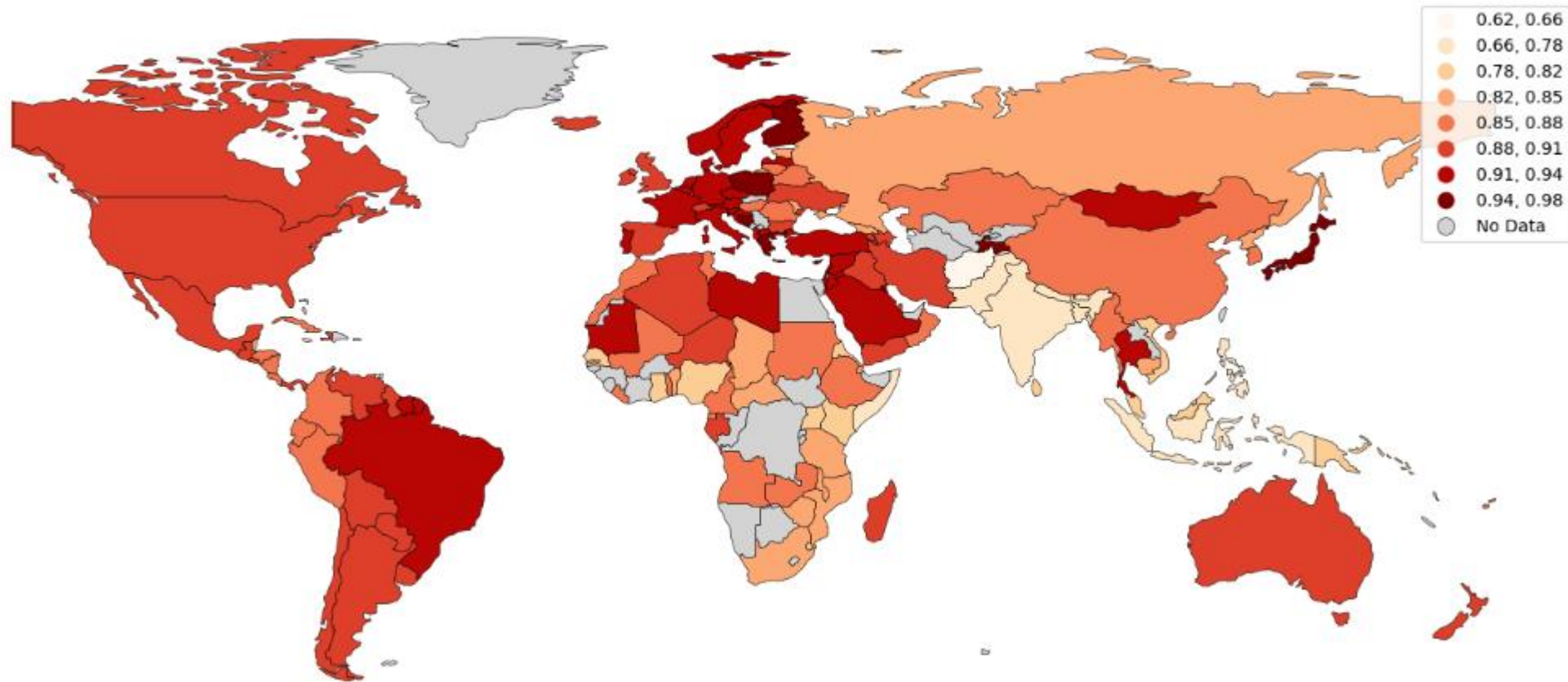
Models = fastText-based classifier

## Are some languages poorly represented?

Model	Lang.	Geo	Non-Geo
Africa, Southern	sot	0.88	0.86
Africa, Sub	bam	0.69	0.53
Africa, Sub	fuh	0.69	0.58
Africa, Sub	ffm	0.70	0.59
Africa, Sub	plt	0.86	0.85
Africa, Sub	tum	0.86	0.83
Africa, Sub	eng	0.87	0.34
Africa, Sub	run	0.89	0.89
America, Central	kek	0.87	0.85
Asia, South	dtv	0.79	0.80
Asia, Southeast	pam	0.81	0.76
Asia, Southeast	cbk	0.85	0.81
Asia, Southeast	spa	0.85	0.41
Asia, Southeast	tet	0.85	0.78
Asia, Southeast	bjn	0.86	0.82
Asia, Southeast	gor	0.87	0.84
Europe, East	eng	0.87	0.34
Europe, East	rmy	0.87	0.75
Europe, Russia	mdf	0.89	0.87
Europe, West	eng	0.86	0.34
Europe, West	ile	0.87	0.88
Oceania	cha	0.86	0.78
Oceania	bjn	0.88	0.82

Table 6: Complete list of all languages with an f-score below 0.90 in the geographic models.

# Where do Geo and Non-Geo Models Disagree?



## Where do Geo and Non-Geo Models Disagree?

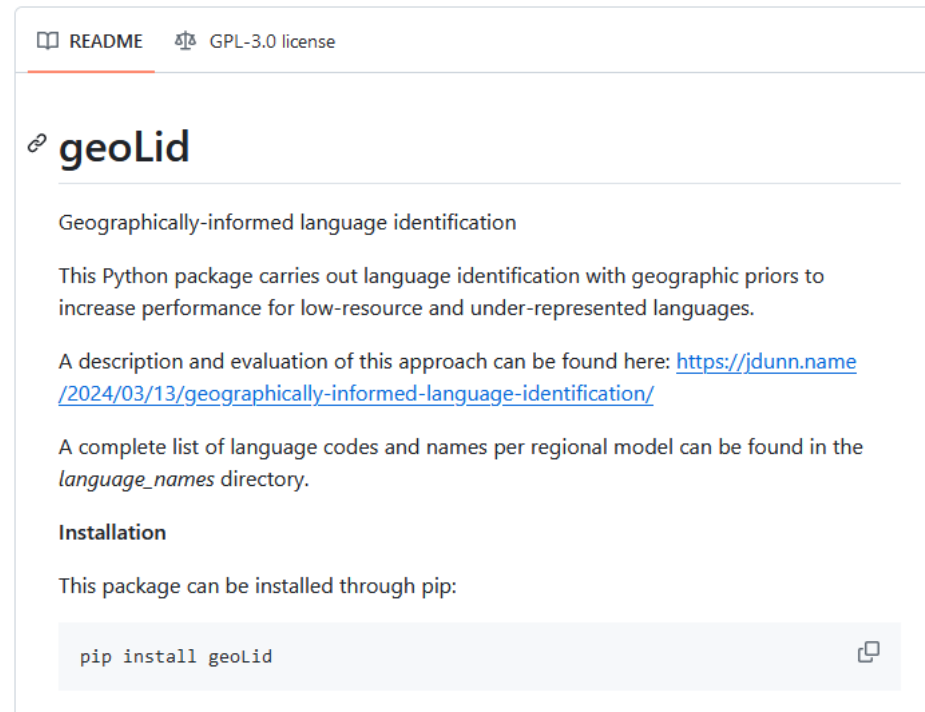
---

<b>Region</b>	<b>N.</b>	<b>Agree</b>	<b>Samples</b>
Africa, North	9	84.93%	10.21 mil
Africa, Southern	3	83.08%	3.77 mil
Africa, Sub	26	84.94%	26.71 mil
America, Central	14	87.28%	16.71 mil
America, North	2	89.62%	2.97 mil
America, South	11	88.85%	14.40 mil
America, Brazil	1	92.41%	1.53 mil
Asia, Central	7	85.52%	9.28 mil
Asia, East	5	89.31%	5.96 mil
Asia, South	7	74.20%	8.42 mil
Asia, Southeast	14	83.86%	13.39 mil
Europe, East	15	89.88%	21.20 mil
Europe, West	23	91.03%	32.19 mil
Europe, Russia,	1	83.03%	1.50 mil
Middle East	12	91.80%	16.33 mil
Oceania	7	85.56%	4.75 mil
<b>Total</b>	<b>157</b>	<b>87%</b>	<b>189 mil</b>

# Conclusions

---

Geographic priors allow better language ID, when available



The screenshot shows the GitHub repository page for 'geoLid'. At the top, there are links for 'README' and 'GPL-3.0 license'. The repository name 'geoLid' is displayed with a link icon. Below the name, the description reads: 'Geographically-informed language identification'. A paragraph follows: 'This Python package carries out language identification with geographic priors to increase performance for low-resource and under-represented languages.' A link is provided: 'A description and evaluation of this approach can be found here: <https://jdunn.name/2024/03/13/geographically-informed-language-identification/>'. Another paragraph states: 'A complete list of language codes and names per regional model can be found in the *language\_names* directory.' The 'Installation' section is titled 'Installation' and contains the text: 'This package can be installed through pip:'. Below this is a code block with the command: `pip install geoLid`. A copy icon is visible to the right of the code block.

[github.com/jonathandunn/geoLid](https://github.com/jonathandunn/geoLid)

**Thanks!**

- Jonathan Dunn

- [jedunn@illinois.edu](mailto:jedunn@illinois.edu)

- [www.jdunn.name](http://www.jdunn.name)

Lane Edwards-Brown

[laneedwardsbrown@gmail.com](mailto:laneedwardsbrown@gmail.com)

