

# Dataset Geography: Mapping Language Data to Language Users

Fahim Faisal, Yinkai Wang, Antonios Anastasopoulos  
Department of Computer Science, George Mason University, USA  
{ffaisal, ywang88, antonis}@gmu.edu

## Abstract

As language technologies become more ubiquitous, there are increasing efforts towards expanding the language diversity and coverage of natural language processing (NLP) systems. Arguably, the most important factor influencing the quality of modern NLP systems is data availability. In this work, we study the geographical representativeness of NLP datasets, aiming to quantify if and by how much do NLP datasets match the expected needs of the language speakers. In doing so, we use entity recognition and linking systems, presenting an approach for good-enough entity linking without entity recognition first. Last, we explore some geographical and economic factors that may explain the observed dataset distributions.<sup>1</sup>

## 1 Introduction

The lack of linguistic, typological, and geographical diversity in NLP research, authorship, and publications is by now widely acknowledged and documented (Caines, 2019; Ponti et al., 2019; Bender, 2011; Adelanı et al., 2021). Nevertheless, the advent of massively multilingual models presents opportunity and hope for the millions of speakers of under-represented languages that are currently under-served by language technologies.

Broadening up the NLP community’s research efforts and scaling from a handful up to the almost 7000 languages of the world is no easy feat. In order for this effort to be efficient and successful, the community needs some necessary foundations to build upon. In seminal work, Joshi et al. (2020) provide a clear overview of where we currently stand with respect to data availability for the world’s languages and relate them to the languages’ representation in NLP conferences. Choudhury and

<sup>1</sup>Code and data are publicly available: [https://github.com/ffaisal93/dataset\\_geography](https://github.com/ffaisal93/dataset_geography). Additional visualizations are available in the project page: <https://nlp.cs.gmu.edu/project/datasetmaps/>.

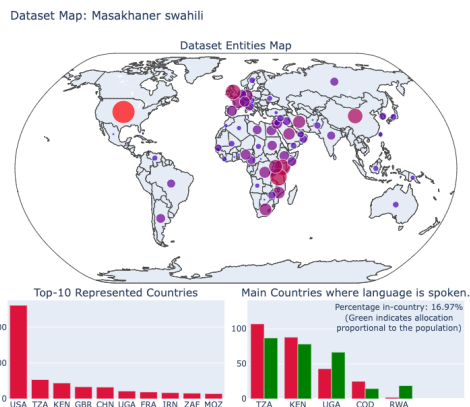


Figure 1: Example of the dataset map our method produces for the Swahili section of MasakhaNER. The dataset is only somewhat representative of Swahili speakers, with only about 17% of entity mentions related to Tanzania, Kenya, Uganda, DR. Congo, or Rwanda and neighboring countries, with the USA and western Europe over-represented.

Deshpande (2021) study how linguistically fair are multilingual language models, and provide a nuanced framework for evaluating multilingual models based on the principles of fairness in economics and social choice theory. Last, Blasi et al. (2022) provide a framework for relating NLP systems’ performance on benchmark datasets to their downstream utility for users at a global scale, which can provide insights into development priorities; they also discuss academic incentives and socioeconomic factors that correlate with the current status of systematic cross-lingual inequalities they observe in language technologies performance.

These works provide insights into current data availability and estimated utility that are paramount for making progress, as well as an evaluation framework for future work. However, there is one missing building block necessary for *real* progress: a way to estimate how representative of the underlying language speakers is the content of our datasets. Any evaluation framework and any utility estimates

we build can only be trustworthy as long as the evaluation data are representative. Gebru et al. (2021) and Bender and Friedman (2018) recognize the importance of this information, including them in their proposed guidelines for “datasheets” and “data statements” respectively; but most datasets unfortunately lack such meta-information. To the best of our knowledge, MaRVL (Liu et al., 2021) is the only dataset that is culturally-aware *by design* in terms of its content.<sup>2</sup>

We propose a method to estimate a dataset’s cultural representativeness by mapping it onto the physical space that language speakers occupy, producing visualizations such as Figure 1. Our contributions are summarized below:

- We present a method to map NLP datasets unto geographical areas (in our case, countries) and use it to evaluate how well the data represent the underlying users of the language. We perform an analysis of the socio-economic correlates of the dataset maps we create. We find that dataset representativeness largely correlates with economic measures (GDP), with geographical proximity and population being secondary.
- We test a simple strategy for performing entity linking by-passing the need for named entity recognition. We evaluate its efficacy on 19 languages, showing that we can get within up to 85% of a NER-informed harder-to-obtain model. We also show that encouragingly, using either model largely leads to similar dataset maps.

## 2 Mapping Datasets to Countries

**Assumptions** This work makes two assumptions: that (a) data locality matters, i.e., speakers of a language are more likely to talk about or refer to local news, events, entities, etc as opposed to ones from a different side of the world, and (b) that we can capture this locality by only focusing on entities. Kumar et al. (2019) discuss these *topical correlations* that are present in datasets,<sup>3</sup> noting that they exist and that L1 language identification models tend to pick up on them, i.e. if a text mentions Finland, a L1 langid model is probably going to predict that the speaker is Finnish, because  $p(\text{Finland}|\text{L1} = \text{Finnish})$  is generally high. In that work Kumar et al. (2019) make explicit effort

<sup>2</sup>Datasets designed to capture dialectal variations, e.g., SD-QA (Faisal et al., 2021), are culturally-aware in terms of annotator selection, but there is no guarantee that their content is also culturally-relevant for the language speakers.

<sup>3</sup>See §2 of their paper.

to avoid learning such correlations because they are interested in building models for  $p(\text{L1}|\text{text})$  (i.e.  $p(\text{L1} = \text{Finnish}|\text{Finland})$ ) that are not confounded by the reverse conditional. The mere fact they need to do this, though, confirms that real-world text has such topical confounds.

As for our second assumption that we can capture these topical correlations by only looking at entities, one need only to take a look at Table 2 of Kumar et al. (2019), which lists the top topical confounding words based on log-odds scores for each L1 language in their dataset: all lists include either entities related to a country where that language is spoken (e.g. ‘Merkel’, the name of a former chancellor, for German) or topical adjectives (e.g. ‘romanian’ for Romanian).

**Approach** For a given dataset, our method follows a simple recipe:

1. Identify named entities present in the dataset.
2. Perform entity linking to wikidata IDs.
3. Use Wikidata to link entities to countries.

We discuss each step below.

**Entity Recognition Step** Standard entity linking is treated as the sequence of two main tasks: entity recognition and entity disambiguation. One approach is to first process the text to extract entities and then disambiguate these entities to the correct entries of a given knowledge base (eg. Wikipedia). This approach relies on NER model quality.

However, to perform analysis on several datasets spanning several low-resource languages, one needs good-quality NER models in all these languages. The interested reader will find a discussion on the cross-lingual consistency of NER models in Appendix F.<sup>4</sup> As we show in Section §4, we can bypass this NER step if we tolerate a small penalty in accuracy.

**Entity Linking Step** In this step we map named entities to their respective Wikidata IDs. We further discuss this step in Section §4.

**From Entities to Countries** We produce maps to visualize the geographical coverage of the datasets we study, discussing their properties and our findings in Section §3.

<sup>4</sup>Discussion summary: state-of-the-art NER models are *not* cross-lingually consistent, i.e. they do not produce the same entity labels when presented with translations of the same sentence. We recommend using parallel data as part of the evaluation sets in multiple languages to measure this important aspect of models’ performance.

To link entities to countries,<sup>5</sup> we rely on Wiki-data entries, depending on the type of entity:

- for persons, we log their places of birth (P19) and death (P20), and country of citizenship (P27);
- for locations, we search for their associated country (P17); and
- for organizations, we use the links of the ‘located\_at’ (P276) and ‘headquartered\_at’ (P159) relations.

Since places of birth/death and headquarters are not necessarily at the country level, we perform a second step of associating these locations with countries. In cases where the result does not correspond to a modern-day country (as can often be the case with historical figures), we do not make any attempts to link it to any modern day countries, excluding them from the analysis.

For example, the entry for Nicolaus Copernicus (Q619) lists him as born in Toruń (Q47554) which is then mapped to Poland; as having died in Frombork (Q497115) that also maps to Poland; and as a citizen of the Kingdom of Poland (Q1649871) which is not mapped to any modern-day country; so he is only linked to Poland. Albert Einstein is similarly mapped to both Germany and the United States, due to his places of birth (Ulm) and death (Princeton).

### 3 Dataset-Country Maps

Before delving into our case studies, we first list a set of statistics of interest that one could extract from our produced dataset-country maps, in order to gauge a dataset’s representativeness.

**Representativeness Measures** We will avoid providing a single metric, largely because the ideal metric to use will be very dataset-specific and related to the goals of the creators of the dataset and the socioeconomic correlates they are interested in (see discussion in Section §3.3).

As a first straightforward representativeness measure, we will compute the **percentage of entities associated with countries where the language is largely spoken**. For example, according to Ethnologue (Eberhard et al., 2021), most Swahili speakers<sup>6</sup> reside in Tanzania, Kenya, Uganda, DR. Congo, and Rwanda. For a Swahili dataset, then, we compute the percentage of all entities associated with this set of countries (“*in-country*”).

<sup>5</sup>A single entity can be associated with a set of more than one countries.

<sup>6</sup>In the case of Swahili they are often second-language speakers.

Notions of equity or fairness across countries could be measured by various fairness metrics, given the distribution of entities over countries in a dataset: from simply computing the standard deviation of the observations,<sup>7</sup> to treating countries as a population and computing fairness indices like the popular Gini index (Gini, 1912; Gastwirth, 1972) or the indices proposed by Speicher et al. (2018). We will opt for a simpler, much more interpretable measure, **the number of countries not represented in the dataset** i.e. countries with associated entity count below a given threshold (we use zero for simplicity but higher values would also be reasonable for large datasets).

Last, especially for languages with significant amounts of speakers in more than one country, it is important to go deeper and measure the representativeness of this *in-country* portion. For a simple example, an English dataset with entities only from the UK is probably not representative of Nigerian or Jamaican English speakers. Hence, we will create two distributions over the countries where the language is largely spoken: the distribution of speaker populations (as available from Ethnologue and other public data), and the distribution of entities observed in the dataset. Discrepancies between these two distributions will reveal potential issues. While one could easily compute some measure of distance between the two distributions (e.g. the Bhattacharyya coefficient (Bhattacharyya, 1943)), in this work we will rely on the interpretable advantages of the visualizations. Measures of fairness could be computed for this portion of the dataset, similarly as discussed above.

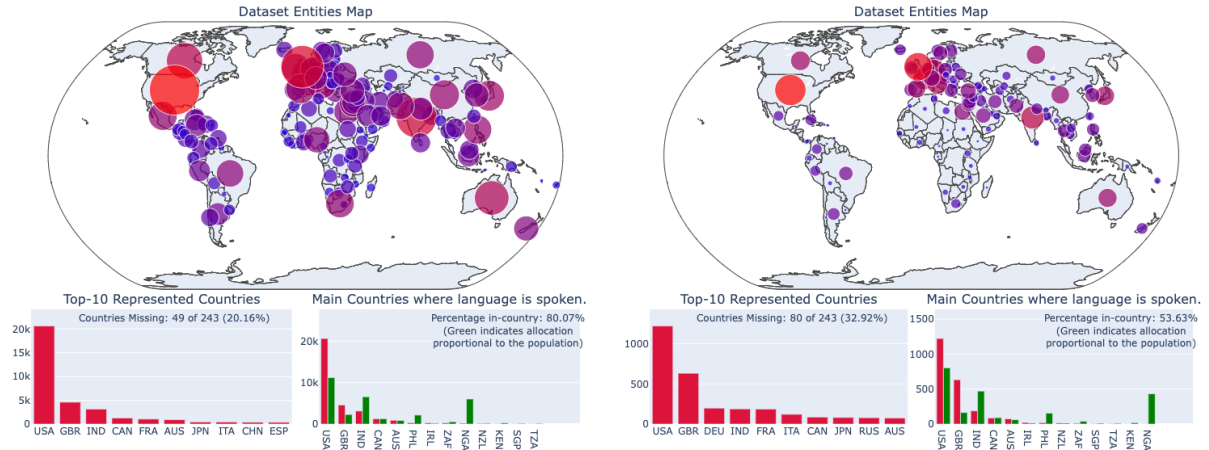
In the example dataset of the Swahili portion of MasakhaNER in Figure 1, the utility of our method is apparent. Through the visualization, a researcher can quickly confirm that the dataset seems to not reflect the users of the language to a large extent: only about 17% of the entities indeed correspond to Tanzania, Kenya, Uganda, DR. Congo, or Rwanda (where Swahili and its varieties are treated as a lingua franca, at least in portions of these countries). Wealthy or populous countries like USA, France, and China, are well-represented,<sup>8</sup> as one would expect, while 156 countries and territories have no representation. At the same time, the visualization allows a researcher to identify gaps:

<sup>7</sup>Or approximations thereof such as the max-min of the observations, as used by (Debnath et al., 2021).

<sup>8</sup>over-represented?

## Natural Questions

## MLQA



## SQuAD

## TyDi-QA (English)

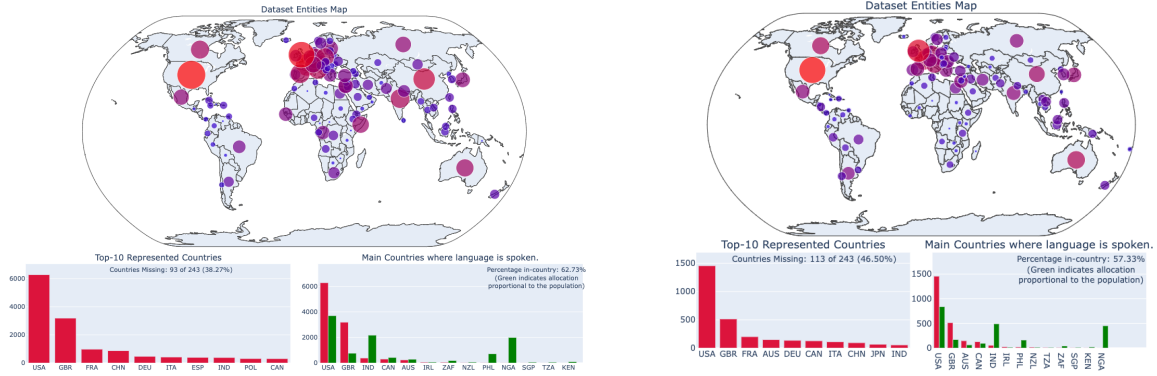


Figure 2: Visualizing the datasets’ geography allows easy comparisons of their representiveness (best viewed in color and zoomed-in). NQ is the most representative of English speakers, with in-country percentage (higher is better) of 80% (SQuAD: 63%; TyDi-QA: 57%; MLQA: 53%) and less countries left unrepresented (lower is better; NQ: 49; MLQA: 80; SQuAD: 93; TyDi-QA: 113).

beyond the neighboring African countries and perhaps the Middle East, north-west African countries as well as central America or central/south-east Asia are clearly under-represented in this dataset. Between the main Swahili-speaking countries, Tanzania, Kenya, and Uganda are well-represented (DR Congo and Rwanda less so, but they have less Swahili speakers), with the former two perhaps slightly over-represented and the latter (as well as Rwanda) being under-represented relative to the speakers population, c.f. red (dataset entities) and green (proportional to population) bars in Figure 1.

### 3.1 Datasets and Settings

We apply the process described above on several datasets, chosen mostly for their language and typological diversity. Our process is not dataset- or

language-dependent,<sup>9</sup> and could easily be applied on any NL dataset. We briefly describe the datasets we include in our study below, with detailed statistics in Appendix C.

**NER Datasets** We study the WikiANN dataset (Pan et al., 2017) that is commonly used in the evaluation of multilingual models. We additionally study the MasakhaNER dataset (Ade-lani et al., 2021), which was created through participatory design (V et al., 2020) in order to focus on African languages. Since these datasets are already annotated with named entities, we only need to perform entity linking.

**Question Answering** We study four question answering datasets (focusing on the questions

<sup>9</sup>Although it does rely on a decent quality entity linker which we lack for most languages. See discussion in §4.

rather than contexts), namely SQuAD (Rajpurkar et al., 2016), MLQA (Lewis et al., 2020), TyDi-QA (Clark et al., 2020), and Natural Questions (Kwiatkowski et al., 2019, NQ;), which have unique characteristics that lend themselves to interesting comparisons. SQuAD is a large English-only dataset (although it has been translated through efforts like XQuAD (Artetxe et al., 2020)). MLQA is a *n*-way parallel multilingual dataset covering 7 languages, created by translating an English dataset. TyDi-QA is another multilingual dataset covering 11 languages, but each language portion is derived separately, without translation involved. Last, NQ is an English QA dataset created based on real-world queries on the Google search engine for which annotators found relevant Wikipedia context, unlike the other datasets that were created by annotators forming questions *given* a context.

**Additional Datasets** While not further discussed in this paper, additional visualizations for more datasets (e.g. for the X-FACTR benchmark (Jiang et al., 2020), and several machine translation benchmarks) are available in the project’s webpage: <https://nlp.cs.gmu.edu/project/datasetmaps/>.

### 3.2 Discussion

Beyond Figure 1, we also show example maps in Figure 2 for NQ, MLQA, SQuAD, and the English portion of TyDi-QA. We provide additional maps for all other datasets in Appendix G.

**Comparing datasets** The comparison of MasakhaNER to the WikiANN dataset (see Appendix G) reveals that the former is rather more localized (e.g. more than 80% of the identified entities in the Dholuo dataset are related to Kenya) while the latter includes a smaller portion from the countries where most native speakers reside (between 10%-20%) and almost always also includes several entries that are very European- or western-centric.

The effect of the participatory design (V et al., 2020) approach on creating the MasakhaNER dataset, where data are curated from local sources, is clear in all language portions of the dataset, with data being highly representative of the speakers. In Figures 8–9 (App. G) the majority of entities in the Wolof portion are from Senegal and neighboring countries (as well as France, the former colonial power of the area), and the Yoruba and Igbo ones are centered on Nigeria.

Figure 2 allows for a direct comparison of different QA datasets (also see maps for other TyDi-QA languages in Appendix G). The first notable point has to do with NQ, which was built based on real-world English-language queries to the Google search engine. Since such queries happen all over the world, this is reflected in the dataset, which includes entities from almost all countries in the world. Two types of countries are particularly represented: ones where English is an official language (USA, UK, Australia, but also, to a lesser extent, India, Nigeria, South Africa, and the Philippines); and wealthy ones (European, Japan, China, etc). In our view, NQ is an exemplar of a representative dataset, because it not only includes representation of most countries where the language is spoken (with the sum of these entities being in their large majority in-country: 80%) but due to its size it also includes entities from almost all countries.

SQuAD also has a large percentage in-country (63%) but it is less representative of different Englishes than NQ. India, for instance, is relatively under-represented in all datasets; in SQuAD it ranks 7<sup>th</sup>, but it ranks 3<sup>rd</sup> in NQ (see red bars in bottom left of figures). On the other hand, the geographical representativeness of both MLQA and TyDi-QA (their English portion) is lacking. Since these datasets rely on Wikipedia articles for their creation, and Wikipedia has a significant western-country bias (Greenstein and Zhu, 2012; Hube and Fetahu, 2018), most entities come from Europe, the US, and the Middle East. All these datasets under-represent English speakers from English-speaking countries of the Global South like Kenya, South Africa, or Nigeria, since there are practically almost no entities from these countries. MLQA further under-represents the speakers of all other languages it includes beyond English, since all data are translations of the English one. Contrast this to TyDi-QA and its visualized Swahili portion which, even though still quite western-centric, does have a higher representation from countries where Swahili is spoken than the TyDi-QA English portion.

This discussion brings forth the importance of being cautious with claims regarding systems’ utility, when evaluated on these datasets. One could argue that a QA system that is evaluated on NQ does indeed give a good estimation of real-world utility; a system evaluated on TyDi-QA gives a distorted notion of utility (biased towards western-based speakers and against speakers from the Global

Factors $\phi$	TyDi-QA (11)		MLQA (1)		SQUAD (1)		NaturalQ. (1)	
	Expl. Var.	MAE	Expl. Var.	MAE	Expl. Var.	MAE	Expl. Var.	MAE
pop	0.272	0.431	0.317	0.401	0.277	1.230	0.395	1.18
gdp	0.507	0.349	0.561	0.332	0.516	1.023	0.535	1.069
gdppc	0.176	0.458	0.182	0.458	0.127	1.345	0.144	1.463
land	0.107	0.504	0.166	0.469	0.142	1.380	0.152	1.459
geo	0.075	0.499	0.040	0.495	0.062	1.393	0.030	1.561
geo+gdp	0.550	0.333	<b>0.579</b>	0.321	<b>0.552</b>	0.932	0.550	1.054
pop+gdp+geo	0.532	0.337	0.548	0.326	0.534	0.940	0.550	1.005
pop+gdp+gdppc+geo	<b>0.555</b>	<b>0.321</b>	0.576	<b>0.310</b>	0.531	<b>0.918</b>	<b>0.570</b>	<b>0.973</b>
all 5 factors	0.538	0.325	0.566	0.312	0.524	0.924	0.561	0.981

Table 1: Empirical comparison of factors on QA datasets, averaging over their respective languages (number in parentheses). We report the five-fold cross-validation explained variance and mean absolute error of a linear model.

South); a system evaluated on MLQA will give an estimation as good as one evaluated on TyDi-QA, but only on the English portion. We clarify that this does not diminish the utility of the datasets themselves as tools for comparing models and making progress in NLP: MLQA is extremely useful for comparing models across languages *on the exact same data*, thus facilitating easy comparisons of the cross-lingual abilities of QA systems, without the need for approximations or additional statistical tests. But we argue that MLQA should not be used to assess the potential utility of QA systems for German or Telugu speakers.

Similar observations can be made about comparing two similar projects that aim at testing the memorization abilities of large language models, namely X-FACTR and multi-LAMA (mLAMA; Kassner et al., 2021) – see corresponding Figures in Appendix G. Both of these build on top of Wikidata and the mTREx dataset. However, mLAMA translates English prompts and uses entity-relation triples mined from the English portion of Wikidata, unlike X-FACTR which uses different data for each language, mined from their respective portion of Wikidata. Both are still western-biased, since they rely on Wikipedia, but one (X-FACTR) is better at giving an indication of potential downstream utility to users.

### 3.3 Socioeconomic Correlates

In this section we attempt to explain our findings from the previous section, tying them to socioeconomic factors.

**Empirical Comparison of Factors** We identify socioeconomic factors  $\phi$  that could be used to explain the observed geographic distribution of the entities in the datasets we study. These are:

- a country’s population  $\phi_{\text{pop}}$
- a country’s gross domestic product (GDP)  $\phi_{\text{gdp}}$
- a country’s GDP per capita  $\phi_{\text{gdppc}}$

- a country’s landmass  $\phi_{\text{land}}$
- a country’s geographical distance from country/ies where the language is spoken  $\phi_{\text{geo}}$

The first four factors are global and fixed. The fifth one is relative to the language of the dataset we are currently studying. For example, when we focus on the Yoruba portion of the mTREx dataset, we use Nigeria (where Yoruba is spoken) as the focal point and compute distances to all other countries. The assumption here is that a Yoruba speaker is more likely to use or be interested in entities first from their home country (Nigeria), then from its neighboring countries (Cameroon, Chad, Niger, Benin) and less likely of distant countries (e.g. Argentina, Canada, or New Zealand). Hence, we assume the probability to be inversely correlated with the country’s distance. For macro-languages or ones used extensively in more than one country, we use a population-weighted combination of the factors of all relevant countries.

To measure the effect of such factors it is common to perform a correlational analysis, where one measures Spearman’s rank correlation coefficient  $\rho$  between the dataset’s observed geographical distribution and the factors  $\phi$ . It is important, though, that the factors are potentially covariate, particularly population and GDP. Hence, we instead compute the variance explained by a linear regression model with factors  $\phi$  as input, i.e.,  $a\phi_{\text{pop}} + b\phi_{\text{gdp}} + c\phi_{\text{gdppc}} + d\phi_{\text{geo}} + e$  with  $a$ – $e$  learned parameters, trained to predict the log of observed entity count of a country. We report explained variance and mean absolute error from five-fold cross-validation experiments to avoid overfitting.

**Socioeconomic Correlates and Discussion** The results with different combination of factors for the QA datasets are listed in Table 1.<sup>10</sup> The best *sin-*

<sup>10</sup>See Appendix H for NER datasets, and Appendix I for a breakdown by language for all datasets.

*gle* predictor is, perhaps unsurprisingly, the GDP of the countries where the language is spoken: all datasets essentially over-represent wealthy countries (e.g. USA, China, or European ones). Note that GDP per capita is not as good a predictor, neither is landmass. A combination of geographical distance with GDP explains most of the variance we observe for all datasets, an observation that confirms the intuitions we discussed before based solely on the visualizations. Importantly, the fact that including population statistics into the model deteriorates its performance is further proof that our datasets are not representative of or proportional to the underlying populations. The only dataset that is indeed better explained by including population (and GDP per capita) is NQ, which we already argued presents an exemplar of representativeness due to its construction protocol.

**Limitations** It is important to note that our assumptions are also limiting factors in our analyses. Mapping languages to countries is inherently lossy. It ignores, for instance, the millions of immigrants scattered throughout the world whose L1 language could be different than the dominant language(s) in the region where they reside. Another issue is that for many languages the necessary granularity level is certainly more fine than country; if a dataset does not include any entities related to the Basque country but does include a lot of entities from Spain and France, our analysis will incorrectly deem it representative, even though the dataset could have been a lot more culturally-relevant for Basque speakers by actually including Basque-related entities.

Another limitation lies in the current state of the methods and data resources on which our approach relies. Beyond discrepancies in NER/EL across languages (addressing which is beyond the scope of this work), we suspect that Wikidata suffers from the same western-centric biases that Wikipedia is known for (Greenstein and Zhu, 2012). As a result, we might be underestimating the cultural representativeness of datasets in low-resource languages.

An additional hurdle, and why we avoid providing a single concrete *representativeness score* or something similar, is that the ideal combination of socioeconomic factors can be subjective. It could be argued, for instance, either that geographic proximity by itself should be enough, or that it should not matter at all. Even further, other factors that we did not consider (e.g. literacy rate or web access) might influence dataset construction

decisions. In any case, we share the coefficients of the NQ model, since it is the most representative dataset we studied, at least for English:  $a = 0.1.46$  (for  $\phi_{pop}$ ),  $b = 0.87$  ( $\phi_{gdp}$ ),  $c = 25.4$  ( $\phi_{gdppc}$ ),  $d = 0.41$  ( $\phi_{geo}$ ). We believe that ideally GDP should not matter ( $b \rightarrow 0$ ) and that a combination of speaker population and geographic proximity is ideal.<sup>11</sup>

### 3.4 Geographical Breakdown of Models' Performance

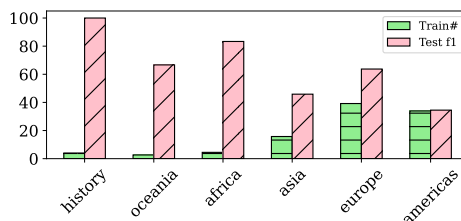
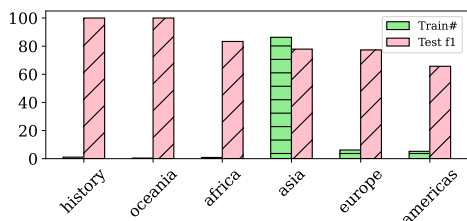
Beyond the analysis of the datasets themselves, we can also break down the performance of models by geographical regions, by associating test (or dev) set samples containing entities with the geographical location of said entities. Since most test sets are rather small (a few hundred to a couple thousand instances) we have to coarsen our analysis: we map each country to a broader region (Africa, Americas, Asia, Europe, Oceania), keeping historical entities in a separate category (History).<sup>12</sup>

We perform such a case study on TyDi-QA, comparing the performance on the TyDi-QA development sets of two models: one trained monolingually on the training set of each language of TyDi-QA (gold task), and another model trained by Debnath et al. (2021) on English SQuAD and automatically generated translations in the target languages. Example results on Telugu shown in Figure 3 reveal some notable trends.<sup>13</sup> First, training set representation (green bars in the Figures) is not a necessary condition for good test set performance (red bars). Some test set instances (e.g. with historical and African entities) receive similar test F1 score from both models. Perhaps the most interesting though, is the comparison of the Asian and European portions of the test set: the Telugu monolingual model achieves similar performance in these two subsets; but the SQuAD-trained model is almost 20 percentage points worse on the Asian subset, showing the potential unfairness of translation-based models (Debnath et al., 2021). For most TyDi-QA languages (Indonesian being an exception, see Table 2) the macro-standard deviation (computed over the averages of the 6 region subsets) is larger for the SQuAD-trained model (which is, hence, less fair than models trained on

<sup>11</sup>However regrettable a fact, it is undeniable that western culture and politics have world-wide effects. So their (over-)representation as a result of their high influence (and GDP) might actually reflect the true interests of people everywhere!

<sup>12</sup>Future work could explore a different clustering.

<sup>13</sup>See Table 4 in Appendix D for all languages.



(a) Train: TyDiQA, Test: TyDiQA dev set (b) Train: (English) SQuAD, Test: TyDiQA dev set

Figure 3: Area-based breakdown of the performance of two models on the Telugu TyDi-QA dev set (red bars) compared with train-set distribution of these geographical areas (green bars). Model (b) is less fair than Model (a). Compare, for instance, the differences in performance between Asia and Europe of the two models.

TyDi-QA Test Set	Stdev over the 6 regions of model trained on		$\Delta$
	SQuAD	TyDi-QA	
Indonesian	17.40	21.52	-4.12
English	13.11	12.66	0.46
Finnish	6.33	5.99	0.3
Arabic	19.24	10.08	9.16
Telugu	21.83	12.45	9.38
Bengali	36.41	10.21	26.1

Table 2: Standard deviation (the lower the more fair the model) of area-based performance averages for two models. Evaluation on TyDi-QA development set.

TyDi-QA).

#### 4 Bypassing NER for Entity Linking

We use mGENRE (Cao et al., 2021) for the task of multilingual entity linking, a sequence to sequence system that predicts entities in an auto-regressive manner. It works particularly well in a zero-shot setting as it considers 100+ target languages as latent variables to marginalize over.

Typically, the input to mGENRE can be informed by a NER model that provides the named entity span over the source. For instance, in the Italian sentence "[START] Einstein [END] era un fisico tedesco." (*Einstein was a German physicist.*) the word Einstein is enclosed within the entity span. mGENRE is trained to use this information to return the most relevant Wikidata entries.

Due to the plasticity of neural models and mGENRE’s auto-regressive token generation fashion, we find that by simply enclosing the whole sentence in a span also yields meaningful results. In particular, for the previously discussed Italian sentence now the input to mGENRE is "[START] Einstein era un fisico tedesco. [END]".

The advantage of this approach is two-fold. First, one does not need a NER component. Second, exactly because of bypassing the NER component, the EL model is now less constrained in its output; in cases where the NER component made errors, there’s a higher chance that the EL model will re-

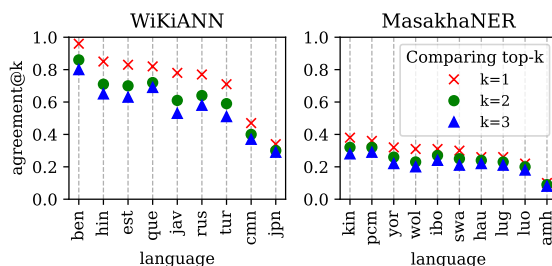


Figure 4: For some languages a NER-Relaxed model is within 60% of a NER-Informed model. agreement@k: ratio of top-k agreement of the models.

turn the correct result.

**Experiments and Results** We conduct experiments to quantify how different a model uninformed by a NER model (NER-Relaxed) will perform compared to one following the typical pipeline (NER-Informed).

Given the outputs of the two models over the same set of sentences, we will compare their average agreement@k, as in the size of the intersection of the outputs of the two models divided by the number of outputs of the NER-Informed model, when focusing only on their top-k outputs.<sup>14</sup> We aggregate these statistics at the sentence level over the whole corpus. We focus on two datasets, namely WikiANN and MasakhaNER, summarizing the results in Figure 4.<sup>15</sup>

Comparing the general performance between these two datasets, it is clear that general agreement is decent. In 7 Out of 9 typologically diverse languages from WikiANN, more than 60% top-1 entities are linked by both models. The African languages from MasakhaNER are low-resource ones yielding less than 40% EL agreement to English in all cases. Given that most of these languages have not been included in the pre-training of BART (the model mGENRE is based on), we expect that using AfriBERTa (Ogueji et al.) or similar models

<sup>14</sup>Both models typically output between 1–3 entity links ranked according to their likelihood.

<sup>15</sup>An extensive results table is available in Appendix E.



in future work would yield improvements.

**Effect on downstream maps** We compare the dataset maps we obtain using `NER-Relaxed` and `NER-Informed` (using gold annotations) models in our pipeline for the `MasakhaNER` dataset. Overall, the maps are very similar. An example visualization of the two maps obtained for Swahili is in Figure 5 in Appendix E.1.

The `NER-Informed` model produces slightly fewer entities overall (likely exhibiting higher precision for lower link recall) but there are minimal differences on the representativeness measures e.g., the in-country percentage changes from 15.3% (`NER-Informed`) to 16.9% (`NER-Relaxed`). We can compare the distributions of the top- $k$  countries obtained with the two models using `Ranked Baised Overlap` (RBO; higher is better; Webber et al., 2010).<sup>16</sup> The results for varying values for  $k$  (top- $k$  countries) are presented in Table 6 in Appendix E.1. We overall obtain very high RBO values ( $> .8$  for  $k = 10$ ) for all language portions and all values of  $k$ . For example for 8 of the 10 `MasakhNER` languages the two models almost completely agree on the top-10 countries with only slight variations in their ranking. `Dholuo` and `Amharic` are the ones exhibiting the worse overlap (but still  $> .5$  RBO).

## 5 Conclusion

We present a recipe for visualizing how representative NLP datasets are with respect to the underlying language speakers. We plan to further improve our tool<sup>17</sup> by making `NER/EL` models more robustly handle low-resource languages. We will also expand our dataset and task coverage, to get a broader overview of the current utility of NLP systems.

## Acknowledgements

This work is generously supported by NSF Awards 2040926 and 2125466.

## References

David Ifeoluwa Adelani, Jade Abbott, Graham Neubig, Daniel D’souza, Julia Kreutzer, Constantine Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, Stephen Mayhew, Israel Abebe Azime, Shamsuddeen Muhammad, Chris Chinenye Emezue, Joyce Nakatumba-Nabende, Perez Ogayo, Anuoluwapo Aremu,

Catherine Gitau, Derguene Mbaye, Jesujoba Alabi, Seid Muhie Yimam, Tajuddeen Gwadabe, Ignatius Ezeani, Rubungo Andre Niyongabo, Jonathan Mukibi, Verrah Otiende, Irero Orife, Davis David, Samba Ngom, Tosin Adewumi, Paul Rayson, Mofetoluwa Adeyemi, Gerald Muriuki, Emmanuel Anebi, Chiamaka Chukwuneke, Nkiruka Odu, Eric Peter Wairagala, Samuel Oyerinde, Clemencia Siro, Tobius Saul Bateesa, Temilola Oloyede, Yvonne Wambui, Victor Akinode, Deborah Nabagereka, Maurice Katusiime, Ayodele Awokoya, Mouhamadane MBOUP, Dibora Gebreyohannes, Henok Tilaye, Kelechi Nwaike, Degaga Wolde, Abdoulaye Faye, Blessing Sibanda, Orevaoghene Ahia, Bonaventure F. P. Dossou, Kelechi Ogueji, Thierno Ibrahima DIOP, Abdoulaye Diallo, Adewale Akinfaderin, Tendai Marengereke, and Salomey Osei. 2021. [Masakhaner: Named entity recognition for african languages](#).

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. [On the cross-lingual transferability of monolingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online. Association for Computational Linguistics.

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. [Findings of the 2020 conference on machine translation \(WMT20\)](#). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55, Online. Association for Computational Linguistics.

Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. [Findings of the 2019 conference on machine translation \(WMT19\)](#). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy. Association for Computational Linguistics.

Rachel Bawden, Giorgio Maria Di Nunzio, Cristian Grozea, Inigo Jauregi Unanue, Antonio Jimeno Yepes, Nancy Mah, David Martinez, Aurélie Névéal, Mariana Neves, Maite Oronoz, Olatz Perez-de Viñaspre, Massimo Piccardi, Roland Roller, Amy Siu, Philippe Thomas, Federica Vezzani, Maika Vicente Navarro, Dina Wiemann, and Lana Yeganova. 2020. [Findings of the WMT 2020 biomedical translation shared task: Basque, Italian and Russian as new additional languages](#). In *Proceedings of the Fifth Conference on Machine Translation*, pages 660–687, Online. Association for Computational Linguistics.

<sup>16</sup>See Appendix E.1 on why this metric is appropriate.

<sup>17</sup>To be released as a Python package.

- Emily M Bender. 2011. On achieving and evaluating language-independence in nlp. *Linguistic Issues in Language Technology*, 6(3):1–26.
- Emily M Bender and Batya Friedman. 2018. Data statements for natural language processing: Toward mitigating system bias and enabling better science. *Transactions of the Association for Computational Linguistics*, 6:587–604.
- A. Bhattacharyya. 1943. On a measure of divergence between two statistical populations defined by their probability distributions. *Bulletin of the Calcutta Mathematical Society*, 35:99–109.
- Federico Bianchi, Debora Nozza, and Dirk Hovy. 2021. [Language invariant properties in natural language processing](#).
- Damián Blasi, Antonios Anastasopoulos, and Graham Neubig. 2022. Systematic inequalities in language technology performance across the world’s languages. In *Proc. ACL*.
- Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleš Tamchyna. 2014. [Findings of the 2014 workshop on statistical machine translation](#). In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. [Findings of the 2017 conference on machine translation \(WMT17\)](#). In *Proceedings of the Second Conference on Machine Translation*, pages 169–214, Copenhagen, Denmark. Association for Computational Linguistics.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. [Findings of the 2016 conference on machine translation](#). In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198, Berlin, Germany. Association for Computational Linguistics.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Barry Haddow, Matthias Huck, Chris Hokamp, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Carolina Scarton, Lucia Specia, and Marco Turchi. 2015. [Findings of the 2015 workshop on statistical machine translation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 1–46, Lisbon, Portugal. Association for Computational Linguistics.
- Ondřej Bojar, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Philipp Koehn, and Christof Monz. 2018. [Findings of the 2018 conference on machine translation \(WMT18\)](#). In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 272–303, Belgium, Brussels. Association for Computational Linguistics.
- Andrew Caines. 2019. [The geographic diversity of nlp conferences](#).
- Nicola De Cao, Ledell Wu, Kashyap Papat, Mikel Artetxe, Naman Goyal, Mikhail Plekhanov, Luke Zettlemoyer, Nicola Cancedda, Sebastian Riedel, and Fabio Petroni. 2021. [Multilingual autoregressive entity linking](#).
- Monojit Choudhury and Amit Deshpande. 2021. How linguistically fair are multilingual pre-trained language models? In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12710–12718.
- Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. TyDi QA: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Arnab Debnath, Navid Rajabi, Fardina Fathmiul Alam, and Antonios Anastasopoulos. 2021. [Towards more equitable question answering systems: How much more data do you need?](#) In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 621–629, Online. Association for Computational Linguistics.
- Zi-Yi Dou and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora. In *Conference of the European Chapter of the Association for Computational Linguistics (EACL)*.
- Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. [A simple, fast, and effective reparameterization of IBM model 2](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.
- David M. Eberhard, Gary F. Simons, and Charles D. (eds.) Fennig. 2021. *Ethnologue, languages of the world*. Twenty-fourth edition. Dallas, Texas: SIL International. Online version: <http://www.ethnologue.com>.

- Fahim Faisal, Sharlina Keshava, Md Mahfuz Ibn Alam, and Antonios Anastasopoulos. 2021. [SD-QA: Spoken dialectal question answering for the real world](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3296–3315, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- ∇, Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohungebe, Solomon Oluwole Akinola, Shamsuddeen Muhammad, Salomon Kabongo Kabenamualu, Salomey Osei, Freshia Sackey, et al. 2020. Participatory research for low-resourced machine translation: A case study in african languages. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 2144–2160.
- Joseph L Gastwirth. 1972. The estimation of the lorenz curve and gini index. *The review of economics and statistics*, pages 306–316.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12):86–92.
- Corrado Gini. 1912. Variabilità e mutabilità. Reprinted in *Memorie di metodologica statistica (Ed. Pizetti E)*.
- Shane Greenstein and Feng Zhu. 2012. Is wikipedia biased? *American Economic Review*, 102(3):343–48.
- Milan Gritta, Mohammad Taher Pilevar, and Nigel Collier. 2019. [A pragmatic guide to geoparsing evaluation: Toponyms, named entity recognition and pragmatics](#). *Language Resources and Evaluation*, 54.
- Matthew Honnibal and Ines Montani. 2017. spacy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.
- Christoph Hube and Besnik Fetahu. 2018. Detecting biased statements in wikipedia. In *Companion proceedings of the the web conference 2018*, pages 1779–1786.
- Zhengbao Jiang, Antonios Anastasopoulos, Jun Araki, Haibo Ding, and Graham Neubig. 2020. [X-FACTR: Multilingual factual knowledge retrieval from pre-trained language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5943–5959, Online. Association for Computational Linguistics.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. [The state and fate of linguistic diversity and inclusion in the NLP world](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.
- Nora Kassner, Philipp Dufter, and Hinrich Schütze. 2021. [Multilingual LAMA: Investigating knowledge in multilingual pretrained language models](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3250–3258, Online. Association for Computational Linguistics.
- Margarita Kokla and Eric Guilbert. 2020. [A review of geospatial semantic information modeling and elicitation approaches](#). *ISPRS International Journal of Geo-Information*, 9(3).
- Sachin Kumar, Shuly Wintner, Noah A. Smith, and Yulia Tsvetkov. 2019. [Topics to avoid: Demoting latent confounds in text classification](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4153–4163, Hong Kong, China. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*.
- Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020. [Mlqa: Evaluating cross-lingual extractive question answering](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7315–7330.
- Fangyu Liu, Emanuele Bugliarello, Edoardo Maria Ponti, Siva Reddy, Nigel Collier, and Desmond Elliott. 2021. [Visually grounded reasoning across languages and cultures](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10467–10485, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Benjamin Muller, Antonios Anastasopoulos, Benoît Sagot, and Djamé Seddah. 2021. [When being unseen from mBERT is just the beginning: Handling new languages with multilingual language models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 448–462, Online. Association for Computational Linguistics.
- Jian Ni, Georgiana Dinu, and Radu Florian. 2017. Weakly supervised cross-lingual named entity recognition via effective annotation and representation projection. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1470–1480.

- Kelechi Ogueji, Yuxin Zhu, and Jimmy Lin. [Small data? no problem! exploring the viability of pretrained multilingual language models for low-resource languages.](#)
- Sebastian Padó and Mirella Lapata. 2009. Cross-lingual annotation projection for semantic roles. *Journal of Artificial Intelligence Research*, 36:307–340.
- Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. [Cross-lingual name tagging and linking for 282 languages.](#) In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1946–1958, Vancouver, Canada. Association for Computational Linguistics.
- Edoardo Maria Ponti, Helen O’Horan, Yevgeni Berzak, Ivan Vulić, Roi Reichart, Thierry Poibeau, Ekaterina Shutova, and Anna Korhonen. 2019. [Modeling language variation and universals: A survey on typological linguistics for natural language processing.](#) *Computational Linguistics*, 45(3):559–601.
- Ross S. Purves, Paul Clough, Christopher B. Jones, Mark H. Hall, and Vanessa Murdock. 2018. [Geographic information retrieval: Progress and challenges in spatial search of text.](#) *Foundations and Trends® in Information Retrieval*, 12(2-3):164–318.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text.](#) In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021. [WikiMatrix: Mining 135M parallel sentences in 1620 language pairs from Wikipedia.](#) In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361, Online. Association for Computational Linguistics.
- Till Speicher, Hoda Heidari, Nina Grgic-Hlaca, Krishna P Gummadi, Adish Singla, Adrian Weller, and Muhammad Bilal Zafar. 2018. [A unified approach to quantifying algorithmic unfairness: Measuring individual & group unfairness via inequality indices.](#) In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2239–2248.
- Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A. Smith, and Yejin Choi. 2020. [Dataset cartography: Mapping and diagnosing datasets with training dynamics.](#) In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9275–9293, Online. Association for Computational Linguistics.
- Jörg Tiedemann. 2014. [Rediscovering annotation projection for cross-lingual parser induction.](#) In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1854–1864.
- William Webber, Alistair Moffat, and Justin Zobel. 2010. [A similarity measure for indefinite rankings.](#) *ACM Transactions on Information Systems (TOIS)*, 28(4):1–38.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Huggingface’s transformers: State-of-the-art natural language processing.](#)

## A Responsible NLP Notes

We use this section to expand on potential limitations and risks of this work.

An inherent limitation of this work is that many datasets are constructed with the goal of answering scientific questions – not necessarily to be used to build NLP systems that serve language users. If our tool is applied without the assumptions behind dataset construction in mind, it might lead to undue criticisms of existing datasets. It is also important to reiterate that no tool, including ours, will ever be 100% accurate, so our tool should be used as *an indicator* of the cultural representativeness of language datasets, not as a tool that can provide definitive answers.

All scientific artifacts used in this paper are publicly available under permissive licenses for fair use. We are not re-distributing any data or code, beyond the code that we wrote ourselves (which will be released under a CC-0 license) and the additional annotations on top of the existing datasets which map the datasets to Wikidata entries (Wikidata data are also available under a CC-0 license). Our use of our data is consistent with their intended use.

## B Related Work

Effective measurement of dataset quality is an aspect of fast-growing significance. Training large language models require huge amount of data and as a result, the inference generated by these pre-trained language model as well as the fine-tuned models often show inherent data bias. In a recent work (Swayamdipta et al., 2020), the authors present how data-quality aware design-decision can improve the overall model performance. They formulated categorization of data-regions based on characteristics such as out-of-distribution feature, class-probability fluctuation and annotation-level discrepancy.

Usually, multilingual datasets are collected from diverse places. So it is important to assess whether the utility of these datasets are representative enough to reflect upon the native speakers. We find the MasakhaNER (Adelani et al., 2021) is one such dataset that was collected from local sources and the data characteristics can be mapped to local users as a result. In addition, language models often requires to be truly language-agnostic depending on the tasks, but one recent work shows that, the current state-of-the-art language applica-

tions are far from achieving this goal (Joshi et al., 2020). The authors present quantitative assessment of available applications and language-resource trajectories which turns out not uniformly distributed over the usefulness of targeted users and speakers from all parts of the world.

Linking dataset entities to geospatial concept is one integral part of our proposed methodology. Ongoing geospatial semantics research mostly focuses on extracting spatial and temporal entities (Kokla and Guilbert, 2020; Purves et al., 2018). The usual approach is to first extract geo-location concepts (i.e. geotagging) from semi-structured as well as unstructured data and then linking those entities to location based knowledge ontology (i.e. geocoding). In (Gritta et al., 2019), the authors propose a task-metric-evaluation framework to evaluate existing NER based geoparsing methods. The primary findings suggest that NER based geo-tagger models in general rely on instant word-sense while avoiding contextual information.

One important aspect of our study is the evaluation of cross-lingual consistency while performing multilingual NER or EI tasks. In (Bianchi et al., 2021), the authors focus on the consistency evaluation of language-invariant properties. In an ideal scenario, the properties should not be changed via the language transformation models but commercially available models are not prone to avoid domain dependency.

## C Dataset Statistics

See details in Table 3.

## D Geographical Breakdown of Models Performance

See details in Table 4.

## E NER-Informed vs NER-Relaxed Models

In this section, we report the detailed results (see Table 5) from our experiment with using intermediate NER model vs skipping this step.

### E.1 Comparison of NER-Informed and NER-Relaxed Maps

This experiment was performed on MasakhaNER data. See Figure 5 for example maps in Swahili. The distributions of the top- $k$  countries we obtain with the two models (one using the gold NER

Dataset	Data-split	Languages	Language count	Sentence count
WikiANN	train	russian, polish, kazakh, bulgarian, finnish, ukrainian, afrikaans, hindi, yoruba, hungarian, dutch-flemish, korean, persian, japanese, javanese, portuguese, hebrew, arabic, spanish-castilian, bengali, urdu, indonesian, tamil, english, malayalam, tagalog, basque, thai, german, romanian-moldavian-moldovan, chinese, telugu, azerbaijani, quechua, modern-greek, turkish, marathi, georgian, estonian, italian, panjabi, burmese, french, gujarati, malay, lithuanian, swahili, vietnamese	48	658600
TyDi-QA	train	english, korean, japanese, telugu, russian, thai, arabic, finnish, bengali, swahili, indonesian	11	166905
MasakhaNER	train	igbo, wolof, nigerian pidgin, kinyarwanda, amharic, hausa, yoruba, ganda, swahili, dholuo	10	12906
SQuAD	train	english	1	130319
MLQA	dev, test	english, simplified chinese, german, arabic, spanish, hindi, vietnamese	7	12738
WMT NEWS	dev, test	polish, kazakh, finnish, xhosa, hindi, japanese, bengali, tamil, zulu, romanian; moldavian; moldovan, chinese, estonian, french, gujarati, inuktitut, lithuanian, turkish, latvian, dholuo, english	20	126972
Natural Questions	train	english	1	307373

Table 3: Statistics of the datasets we study.

	europa	asia	africa	americas	history	oceania
swahili	(80.3, 88.9)	(64.1, 83.4)	(75.5, 81.4)	(88.1, 89.3)	(83.3, 100)	(86.5, 81.2)
bengali	(60.0, 79.6)	(71.0, 79.5)	-	-	(100, 100)	(0, 100)
arabic	(65.2, 79.0)	(74.5, 82.6)	(72.3, 79.0)	(82.4, 82.6)	(36.3, 65.6)	(100, 100)
korean	(19.3, 23)	(30.4, 36.5)	(0, 0)	(23.9, 24.6)	(42.9, 52.4)	-
english	(74.7, 89.2)	(84.0, 80.2)	(60.0, 60.0)	(75.6, 82.9)	(100, 100)	(93.3, 93.3)
indonesian	(79.4, 88.5)	(75.3, 84)	(80, 100)	(79.9, 84.7)	(83.3, 66.7)	(33.3, 33.3)
russian	(65.1, 80.1)	(59.6, 79.1)	(64.9, 67.8)	(67.8, 81.8)	(47.2, 72.3)	(76.8, 66.7)
telugu	(63.7, 77.3)	(45.9, 77.9)	(83.3, 83.3)	(34.5, 65.7)	(100, 100)	(66.7, 100)
finnish	(73.4, 81)	(86.2, 88.9)	(81, 91.7)	(75.9, 83)	(67.7, 74.7)	-

Table 4: Detailed Breakdown of area-based performance (f1 score) of two trained QA models (TyDi-QA, SQuAD). Evaluation is performed on TyDi-QA development set (gold task).

annotations for NEL and one using our NER-relaxed approach) are compared using Ranked Biased Overlap (RBO; higher is better) (Webber et al., 2010), a metric appropriate for computing the weighted similarity of disjoint rankings. We choose a “weighted” metric because we care more about having similar results in the top- $k$  countries (the ones most represented) so that the metric is not dominated by the long tail of countries that may have minimal representation and thus similar rank. We also need a metric that can handle disjoint rankings, since there’s no guarantee that the top- $k$  countries produced by the processes using different

models will be different.<sup>18</sup>

The results for varying values for  $k$  (top- $k$  countries) are presented in Table 6. We overall obtain very high RBO values ( $> .75$ ) for all language portions and all settings.

## F On the Cross-Lingual Consistency of NER/EL Models

**Definition** Bianchi et al. (2021) in concurrent work point out the need to focus on consistency evaluation of **language-invariant properties (LIP)**: properties which should not be changed via language transformation models. They suggest

<sup>18</sup>Metrics like Kendall’s  $\tau$  would suffer from both issues.

Language	k=1	k=2	k=3	Dataset
hin	(4239, 761, 0.85)	(6765, 2717, 0.71)	(8377, 4436, 0.65)	WikiANN
cmn	(9354, 10646, 0.47)	(16015, 23899, 0.4)	(21835, 37346, 0.37)	
jpn	(6739, 13259, 0.34)	(12148, 27820, 0.3)	(17220, 42463, 0.29)	
rus	(15325, 4675, 0.77)	(24663, 13989, 0.64)	(31520, 23051, 0.58)	
est	(16687, 3313, 0.83)	(24413, 10536, 0.7)	(28146, 16459, 0.63)	
ben	(9575, 425, 0.96)	(15759, 2541, 0.86)	(20106, 4930, 0.8)	
que	(82, 18, 0.82)	(124, 48, 0.72)	(159, 72, 0.69)	
tur	(14206, 5794, 0.71)	(21165, 14999, 0.59)	(25053, 23597, 0.51)	
jav	(78, 22, 0.78)	(103, 67, 0.61)	(113, 101, 0.53)	
pcm	(549, 994, 0.36)	(955, 2033, 0.32)	(1217, 3030, 0.29)	
kin	(593, 952, 0.38)	(924, 1988, 0.32)	(1112, 2853, 0.28)	
wol	(242, 534, 0.31)	(350, 1158, 0.23)	(435, 1692, 0.2)	
hau	(417, 1178, 0.26)	(747, 2333, 0.24)	(941, 3402, 0.22)	
ibo	(494, 1093, 0.31)	(834, 2225, 0.27)	(1056, 3257, 0.24)	
amh	(117, 1088, 0.1)	(210, 2184, 0.09)	(289, 3198, 0.08)	
swa	(499, 1175, 0.3)	(819, 2445, 0.25)	(1007, 3678, 0.21)	
lug	(283, 824, 0.26)	(486, 1657, 0.23)	(644, 2362, 0.21)	
yor	(430, 894, 0.32)	(673, 1909, 0.26)	(839, 2893, 0.22)	
luo	(122, 428, 0.22)	(207, 844, 0.2)	(264, 1184, 0.18)	

Table 5: Breakdown of entity extraction count while using NER-informed model. Here for each top k extracted entities, the triplet is the aggregated value of (count of common entities extracted by both ner-informed and ner-relaxed models, count of entities only extracted by ner-relaxed models, ratio of common entity count and total top-k extract by ner-relaxed model )

LIPs include meaning, topic, sentiment, speaker demographics, and logical entailment We propose a definition tailored to entity-related tasks: cross-lingual consistency is the desirable property that two parallel sentences in two languages, which should in principle use the same named entities (since they are translations of each other), are actually tagged with the same named entities.

## F.1 NER Experiments

**Models** We study two models: SpaCy (Honnibal and Montani, 2017): a state-of-art monolingual library that supports several core NLP tasks; and a mBERT-based NER model trained on datasets from WikiANN using the transformers library (Wolf et al., 2020).

**Training** To task-tune the mBERT-based model on the NER task we use the WikiANN dataset with data from the four languages we study: Greek (el), Italian (it), Chinese (zh), and English (en).

**Evaluation** To evaluate cross-lingual consistency, ideally one would use parallel data where both sides are annotated with named entities. What we use instead, since such datasets do not exist to the best of our knowledge, is ‘silver’ annotations over parallel data. We start with unannotated parallel data from the WikiMatrix dataset (Schwenk et al., 2021) and we perform NER on both the English and the other language side, using the respective language model for each side.

In the process of running our experiments, we identified some sources of noise in the WikiMatrix dataset (e.g. mismatched sentences that are clearly not translations of each other). Thus, we calculated the average length ratio between two matched sentences, and discarded data that diverged by more than one standard deviation from the mean ratio, in order to keep 95% of the original data that are more likely to indeed be translations of each other.

We use the state-of-the-art AWESOME-align tool (Dou and Neubig, 2021) as well fast-align (Dyer et al., 2013) to create word-level links between the words of each English sentence to their corresponding translations. Using these alignment links for cross-lingual projection (Padó and Lapata, 2009; Tiedemann, 2014; Ni et al., 2017, *inter alia*) allows us to calculate cross-lingual consistency, measuring the portion of labels that agree following projection. In particular, we use the cross-lingual projections from the English side as ‘correct’ and measure precision, recall, and F-score against them.

**Results** In preliminary experiments we found that, consistently with the literature, AWESOME-align performed generally better than fast-align, hence for the remainder of our experiments we only use AWESOME-align.

For the three languages we study, the cross-lingual consistency of the monolingual SpaCy models is really low, with scores of 8.6% for Greek–

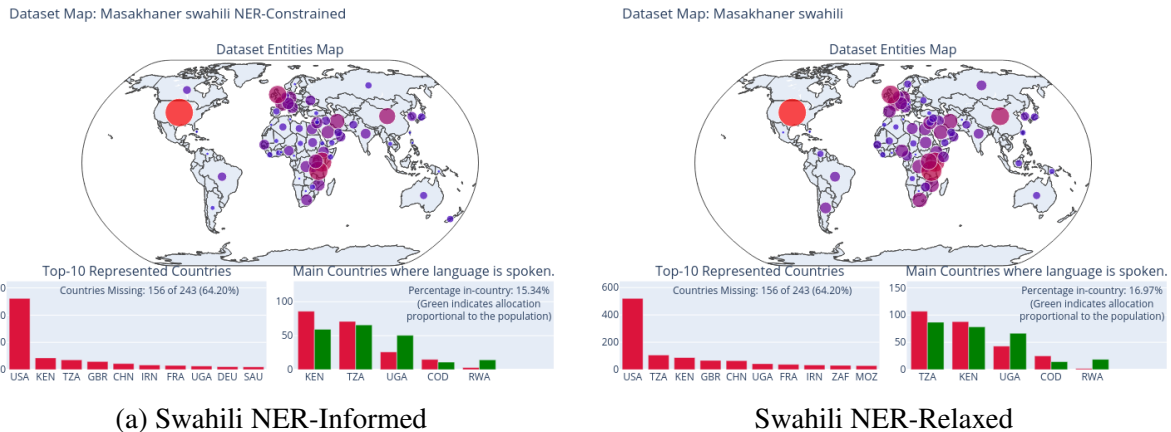


Figure 5: The dataset maps obtained by NER-Informed and NER-Relaxed are very similar, with very small differences in the representativeness measures.

Dataset Portion	Rank Biased Overlap (RBO) for top- $k$ ranked countries with $k=$									
	1	2	3	5	10	20	50	100	200	
Amharic	0.00	0.25	0.50	0.57	0.53	0.51	0.59	0.65	0.76	
Yoruba	1.00	1.00	1.00	0.87	0.82	0.87	0.85	0.83	0.87	
Hausa	1.00	1.00	1.00	0.87	0.80	0.80	0.83	0.83	0.88	
Igbo	1.00	1.00	1.00	0.96	0.89	0.82	0.79	0.79	0.86	
Kinyarwanda	1.00	0.75	0.83	0.86	0.91	0.89	0.83	0.80	0.86	
Luganda	1.00	0.75	0.83	0.81	0.81	0.82	0.78	0.76	0.83	
Dholuo	1.00	0.75	0.83	0.77	0.66	0.58	0.57	0.62	0.76	
Nigerian Pidgin	1.00	1.00	1.00	0.95	0.91	0.90	0.89	0.86	0.90	
Wolof	1.00	1.00	1.00	0.96	0.85	0.77	0.68	0.70	0.81	
Swahili	1.00	0.75	0.83	0.90	0.89	0.89	0.84	0.85	0.90	
Average	0.90	0.82	0.88	0.85	0.81	0.79	0.76	0.77	0.84	

Table 6: Rank Biased Overlap (RBO; higher is better) for the top- $k$  ranked countries obtained by a NER-Informed and a NER-Relaxed model on the MasakhaNER datasets.

Model	Greek	Italian	Chinese
Monolingual (SpaCy)	8.6	3.1	14.1
mBERT	<b>53.4</b>	<b>62.9</b>	<b>25.5</b>

Table 7: Using a multilingual NER model leads to significantly higher consistency tested on Eng-X data.

English, 3.1% for Italian-English and 14.1% for Chinese-English. The SpaCy models are independently trained for each language and can produce 18 fine-grained NE labels e.g. distinguishing dates from time, or locations to geopolitical entities. As such, there was no a priori expectation for high cross-lingual consistency. Nevertheless, these extremely low scores reveal deeper differences, such as potentially widely different annotation protocols

across languages.<sup>19</sup>

For the mBERT-based model we again label both sides of the parallel data, but now evaluate only on locations (LOC), organizations (ORG) and persons (PER) (the label types present in WikiANN). The mBERT models have significantly higher cross-lingual consistency: on the same dataset as above, we obtain 53.4% for Greek to English, 62.9% for Italian to English and 25.5% for Chinese to English.

**Discussion** To further understand the source of cross-lingual discrepancies, we performed manual analysis of 400 Greek-English parallel sentences where the mBERT-based model’s outputs on Greek and the projected labels through English

<sup>19</sup>We note that our evaluation does focus only on labels shared between models/languages.



disagreed.<sup>20</sup> We sampled 100 sentences where the English-projected label was  $\emptyset$  but the Greek one was LOC (location), 100 sentences with English-projected as LOC but Greek as  $\emptyset$ , and similarly for persons (PER).

We performed annotation using the following schema:

- Greek wrong: for cases where only the English-side projected labels are correct
- English wrong: for cases where the English-side projected labels are wrong but the Greek-side are correct
- both wrong: for cases where the labels on both sides are incorrect
- alignment wrong: for cases where the two aligned phrases are not translations of each other, so we should not take the projected labels into account nor compare against them.
- all correct: both sides as well as the alignments are correctly tagged (false negatives).

Encouragingly, the entity alignments were wrong in less than 10% of the parallel sentences we manually labelled. This means that our results are quite robust: a 10%-level of noise cannot account for an almost 50% lack of consistency on the Greek-English dataset.<sup>21</sup> Hence, the system definitely has room for improvement. A second encouraging sign is that less than 2% of the cases were in fact false negatives, i.e. due to the phrasing of the translation only one of the two sides actually contained an entity.

Going further, we find that mistakes vary significantly by label type. In about 75% of the  $\emptyset$ -LOC cases it was the Greek-side labels that were wrong in outputting LOC tags. A common pattern (about 35% of these cases) was the Greek model tagging months as locations. In the case of  $\emptyset$ -PER cases, 62% of the errors were on the English side. A common pattern was the English-side model not tagging persons when they are the very first token in a sentence, i.e. the first token in ‘Olga and her husband [...]’. Appendix K extends this discussion with additional details and examples.

The above observations provide insights into NER models’ mistakes, which we were able to easily identify by contrasting the models’ predictions over parallel sentences. We argue this proves the utility and importance of also evaluating NER mod-

<sup>20</sup>We chose this language pair because one of the authors is a fluent speaker of both languages.

<sup>21</sup>It does provide a potential upper bound of around 90% on the consistency we should expect to find.

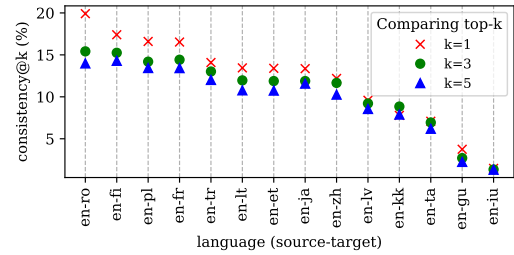


Figure 6: The entity linking cross-lingual consistency is generally low across languages, but especially for low-resource language pairs like English to Inuktitut (iu), Gujarati (gu), or Tamil (ta).

els against parallel data even without gold NER annotations. Improving the NER cross-lingual consistency should in principle also lead to better NER models in general. Potential solutions could use a post-pretraining alignment-based fine-tuned mBERT model as the encoder for our data, or operationalize our measure of cross-lingual consistency into an objective function to optimize.<sup>22</sup>

## F.2 Entity Linking Experiments

We now turn to entity linking (EL), evaluating mGENRE’s cross-lingual consistency (under the NER-Relaxed setting, so the results below should be interpreted under this lens, as the NER-Informed—which we cannot run due to the lack of NER models for some languages—could very well yield different results and analysis).

**Dataset** We use parallel corpora from the WMT news translation shared tasks for the years 2014 to 2020 (Bojar et al., 2014, 2015, 2016, 2017, 2018; Barrault et al., 2019, 2020). We work with 14 English-to-target language pairs, with parallel sentence counts in the range of around 1-5k.

**Evaluation** Unlike our NER experiment settings, we do not need word-level alignments to calculate cross-lingual consistency. We can instead compare the sets of the linked entities for both source and target sentences. As before, we use mGENRE in a NER-Relaxed manner. In an ideal scenario, the output of the model over both source and target language sentences will include the same entity links, yielding a perfect cross-lingual consistency score of 1. In this manner, we calculate and aggregate sentence-level scores for the top- $k$  linked entities for  $k = 1, 3, 5$ . In Figure 6, we present this score as a percentage, dividing the size of the intersection

<sup>22</sup>We leave this for future work, as it detracts off the main goal of this work (mapping datasets to the language users and measuring their representativeness).

src-tgt	k=1 %	k=3 %	k=5 %	sentence count
en-ro	19.91	15.42	13.98	1999
en-fi	17.40	15.25	14.29	1500
en-pl	16.60	14.19	13.43	2000
en-fr	16.53	14.42	13.42	1500
en-tr	14.09	13.02	12.01	1001
en-lt	13.45	11.96	10.77	2000
en-et	13.40	11.88	10.74	2000
en-ja	13.36	11.88	11.57	1998
en-zh	12.19	11.66	10.26	2002
en-lv	9.59	9.21	8.55	2003
en-kk	7.79	8.84	7.88	2066
en-ta	7.09	6.94	6.19	1989
en-gu	3.75	2.70	2.24	1998
en-iu	1.47	1.34	1.31	5173

Table 8: Cross-lingual consistency score (%) for top-k extracted and linked entities over all source language sentences.

(of the source and target sentence outputs) by the number of source sentence entities.

Additionally, in Table 8, we report the detailed cross-lingual consistency score percentages for 14 english-language source-target pairs from WMT news translation shared tasks (Bawden et al., 2020).

**Results** As Figure 6 shows, we obtain low consistency scores across all 14 language pairs, ranging from 19.91% for English-Romanian to as low as 1.47% for English-Inuktitut ( $k = 1$ ). The particularly low scores for languages like Inuktitut, Gujarati, and Tamil may reflect the general low quality of mGENRE for such languages, especially because they use non-Latin scripts, an issue already noted in the literature (Muller et al., 2021).

The low percentage consistency scores for all languages makes it clear that mGENRE does not produce similar entity links for entities appearing in different languages. In future work, we plan to address this limitation, potentially by weighting linked-entities according to the cross-lingual consistency score when performing entity disambiguation in a multilingual setting.

**Discussion** We further analyze whether specific types of entities are consistently recognized and linked across language. We use SpaCy’s English NER model to categorize all entities. Figure 7 presents a visualization comparing consistent entity category counts to source-only ones.

Entity category	Common	Source-only
Unknown	1720	16709
PERSON	1358	5713
ORG	1047	6911
GPE	666	7379
NORP	176	1895
DATE	102	1427
CARDINAL	78	565
EVENT	77	777
LOC	62	453
WORK_OF_ART	20	133
PRODUCT	15	91
FAC	14	161
QUANTITY	8	85
TIME	6	43
MONEY	4	14
LAW	3	113
LANGUAGE	3	80
ORDINAL	2	90
PERCENT	1	3
TOTAL	5362	42642

Table 9: SpaCy NER (Honnibal and Montani, 2017) defined types and counts for consistent linked entities.

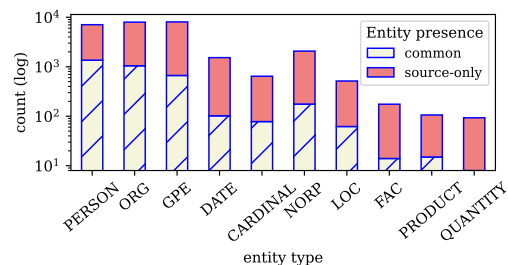


Figure 7: Counts of linked entity types across all WMT language pairs. Notice the y-axis log-scale: many entities are linked differently on non-English input.

From Figure 7, it is clear that geopolitical entities (GPE) are the ones suffering the most from low cross-lingual consistency, with an order of magnitude less entities linked on both the English and the other language side. On the other hand, person names (PER) seem to be easier to link. While the most common types of entities are PERSON, ORG (i.e. organization) and GPE (i.e. geopolitical entity), we found that the NER model still failed to correctly categorize entities like (Surat, Q4629, LOC), (Aurangzeb, Q485547, PER). However, these entities were correctly linked by the NER-Relaxed pipeline, indicating its usefulness. We hypothesize, and plan to test in future work, that a NER-Relaxed entity

further regularized towards cross-lingual consistency will perform better than a NER-Informed pipeline, unless the NER component also shows improved cross-lingual consistency.

From Figure 7, it is clear that geopolitical entities (GPE) are the ones suffering the most from low cross-lingual consistency, with an order of magnitude less entities linked on both the English and the other language side. On the other hand, person names (PER) seem to be easier to link. While the most common types of entities are PERSON, ORG (i.e. organization) and GPE (i.e. geopolitical entity), we found that the NER model still failed to correctly categorize entities like (Surat, Q4629, LOC), (Aurangzeb, Q485547, PER). However, these entities were correctly linked by the NER-Relaxed pipeline, indicating its usefulness. We hypothesize, and plan to test in future work, that a NER-Relaxed entity further regularized towards cross-lingual consistency will perform better than a NER-Informed pipeline, unless the NER component also shows improved cross-lingual consistency.

## G Additional Dataset Maps

We present all dataset maps for the datasets we study:

- MasakhaNER languages are available in Figures 8 and 9.
- TydiQA languages are available in Figures 10 and 11.
- WikiANN (panx) languages are available in Figures 12 through 16.
- SQuAD (English) in Figure 17.

## H NER Dataset Socioeconomic Factors

Table 1 presents the same analysis as the one described in Section 3.3 for the X-FACTR and the NER datasets. The trends are similar to the QA datasets, with GDP being the best predictor and including population statistics hurting the explained variance.

## I Socioeconomic Correlates Breakdown

You can find the breakdown of the socioeconomic correlates in Table 12 for TyDi-QA, Table 13 for MasakhaNER, and Table 14 for WikiANN.

## J NER Models Confusion Matrices

See Figure 18 for the confusion matrices of the SpaCy and our WikiANN neural model.

## K Greek-English NER Error Discussion

We find that the mistakes we identify vary significantly by label. In about 75% of the  $\theta$ -LOC cases it was the Greek-side labels that were wrong in tagging a span as a location. A common pattern we identified (about 35% of these cases) was the Greek model tagging as location what was actually a month. For instance, in the sentence *Ton Máio tu 1990 episkéftikan yia tésseris iméres tin Ouggaria* (*In May 1990, they visited Hungary for four days.*) the model tags the first two words (“in May”) as a location, while the English one correctly leaves them unlabelled.

In the case of LOC- $\theta$  cases, we found an even split between the English- and the Greek-side labels being wrong (with about 40% of the sentences each). Common patterns of mistakes in the English side include tagging persons as locations (e.g. “Heath” in “Heath asked the British to heat only one room in their houses over the winter.” where “Heath” corresponds to Ted Heath, a British politician), as well as tagging adjectives, often locative, as locations, such as “palaeotropical” in “Palaeotropical refers to geographical occurrence.” and “French” in “A further link [...] by vast French investments and loans [...]”.

Last, in the case of  $\theta$ -PER cases we studied, we found that 62% of the errors were on the English side. A common pattern was the English-side model not tagging persons when they are the very first token in a sentence, i.e. the first tokens in “Olga and her husband were left at Ay-Todor.”, in “Friedman once said, ‘If you want to see capitalism in action, go to Hong Kong.’”, and in “Evans was a political activist before [...]” were all tagged as  $\theta$ . To a lesser extent, we observed a similar issue when the person’s name followed punctuation, e.g. “Yavlinsky” in the sentence “In March 2017, Yavlinsky stated that he will [...]”.

## L Comparing X-FACTR to mLAMA

These two similar projects aim at testing the memorization abilities of large language models (X-FACTR and multi-LAMA (mLAMA; Kassner et al., 2021)) – see corresponding Figures in Table ???. Both of these build on top of Wikidata and the mTREx dataset. Hence, their English portions are equally representative of English speakers, suffering from under-representation of English speakers of the Global South. For the other language, however, mLAMA translates English

Factors $\phi$	X-FACTR (11)		MasakhaNER (10)		WikiANN (48)	
	Explained Variance	MAE	Explained Variance	MAE	Explained Variance	MAE
pop	0.356	0.457	0.300	0.295	0.387	0.470
gdp	0.516	0.407	0.341	0.295	0.575	0.382
geo	0.022	0.585	0.100	0.359	0.069	0.586
pop+gdp	0.495	0.403	0.348	0.285	0.553	0.388
pop+geo	0.356	0.455	0.369	0.290	0.399	0.467
geo+gdp	<b>0.521</b>	<b>0.398</b>	<b>0.443</b>	<b>0.284</b>	<b>0.591</b>	<b>0.376</b>
pop+gdp+geo	0.504	<b>0.398</b>	0.440	0.285	0.572	0.380

Table 10: Empirical comparison of factors on NER datasets, averaging over their respective languages (number in parentheses). We report the five-fold cross-validation explained variance and mean absolute error of a linear model.

		geo+gdp	
Language	Country	Expl. Var.	Mean Error
Greek	GRC	0.586	0.343
Yoruba	NGA	0.575	0.219
Bengali	BGD	0.552	0.349
Marathi	IND	0.587	0.29
French	FRA	0.569	0.452
Hebrew	ISR	0.604	0.369
Hungarian	HUN	0.621	0.375
Russian	RUS	0.601	0.406
Spanish	ESP	0.552	0.457
Turkish	TUR	0.613	0.36
Vietnamese	VNM	0.521	0.398
Average		0.504	0.398

Table 11: Language breakdown of the most predictive factors ( $\phi_{\text{geo}}$  and  $\phi_{\text{gdp}}$ ) on X-FACTR dataset.

prompts and uses entity-relation triples mined from the English portion of Wikidata, unlike X-FACTR which uses different data for each language, mined from their respective portion of Wikidata. Both are still western-biased, since they rely on Wikipedia, but one (X-FACTR) is better at giving an indication of potential downstream utility to users.

		geo+gdp	
Language	Country	Expl. Var.	Mean Error
Arabic	SAU	0.501	0.415
Bengali	BGD	0.498	0.385
English	USA	0.562	0.335
Finnish	FIN	0.566	0.376
Indonesian	IDN	0.515	0.387
Japanese	JPN	0.558	0.388
Korean	KOR	0.546	0.336
Russian	RUS	0.522	0.400
Swahili	KEN	0.428	0.469
Telugu	IND	0.534	0.294
Thai	THA	0.550	0.333
Average		0.550	0.333

Table 12: Language breakdown of the most predictive factors ( $\phi_{\text{geo}}$  and  $\phi_{\text{gdp}}$ ) on the TyDi-QA dataset.

		geo+gdp	
Language	Country	Expl. Var.	Mean Error
Amharic	ETH	0.131	0.220
Yoruba	NGA	0.338	0.258
Hausa	NGA	0.321	0.317
Igbo	NGA	0.326	0.207
Kinyarwanda	RWA	0.198	0.229
Luganda	UGA	0.302	0.195
Luo	ETH	0.000	0.110
Nigerian English	NGA	0.493	0.231
Wolof	CMR	0.378	0.160
Swahili	KEN	0.443	-0.285
Average		0.378	0.160

Table 13: Language breakdown of the most predictive factors ( $\phi_{\text{geo}}$  and  $\phi_{\text{gdp}}$ ) on MasakhaNER dataset.

geo+gdp			
Language	Country	Expl. Var.	Mean Error
af	ZAF	0.497	0.338
ar	SAU	0.570	0.454
az	AZE	0.566	0.395
bg	BGR	0.511	0.475
bn	BGD	0.442	0.502
de	DEU	0.613	0.402
el	GRC	0.484	0.456
es	ESP	0.497	0.462
et	EST	0.565	0.398
eu	ESP	0.565	0.387
fa	IRN	0.589	0.426
fi	FIN	0.590	0.411
fr	FRA	0.597	0.408
gu	IND	0.068	0.030
he	ISR	0.551	0.456
hi	IND	0.529	0.279
hu	HUN	0.563	0.451
id	IDN	0.488	0.442
it	ITA	0.569	0.436
ja	IDN	0.591	0.343
jv	JPN	0.062	0.069
ka	GEO	0.474	0.435
kk	KAZ	0.411	0.205
ko	KOR	0.519	0.423
lt	LTU	0.533	0.395
ml	IND	0.495	0.367
mr	IND	0.530	0.320
ms	MYS	0.496	0.463
my	MMR	0.105	0.038
nl	NLD	0.582	0.435
pa	IND	0.052	0.064
pl	POL	0.584	0.436
pt	PRT	0.567	0.432
qu	PER	0.301	0.090
ro	ROU	0.581	0.436
ru	RUS	0.576	0.435
sw	KEN	0.402	0.223
ta	LKA	0.524	0.367
te	IND	0.351	0.107
th	THA	0.567	0.215
tl	PHL	0.473	0.399
tr	TUR	0.619	0.409
uk	UKR	0.576	0.447
ur	PAK	0.512	0.463
vi	VNM	0.557	0.440
yo	NGA	0.079	0.086
zh	CHN	0.591	0.376
Average		0.591	0.376

Table 14: Language breakdown of the most predictive factors ( $\phi_{\text{geo}}$  and  $\phi_{\text{gdp}}$ ) on the WikiANN dataset. 3401

## MasakhaNER Geographic Coverage

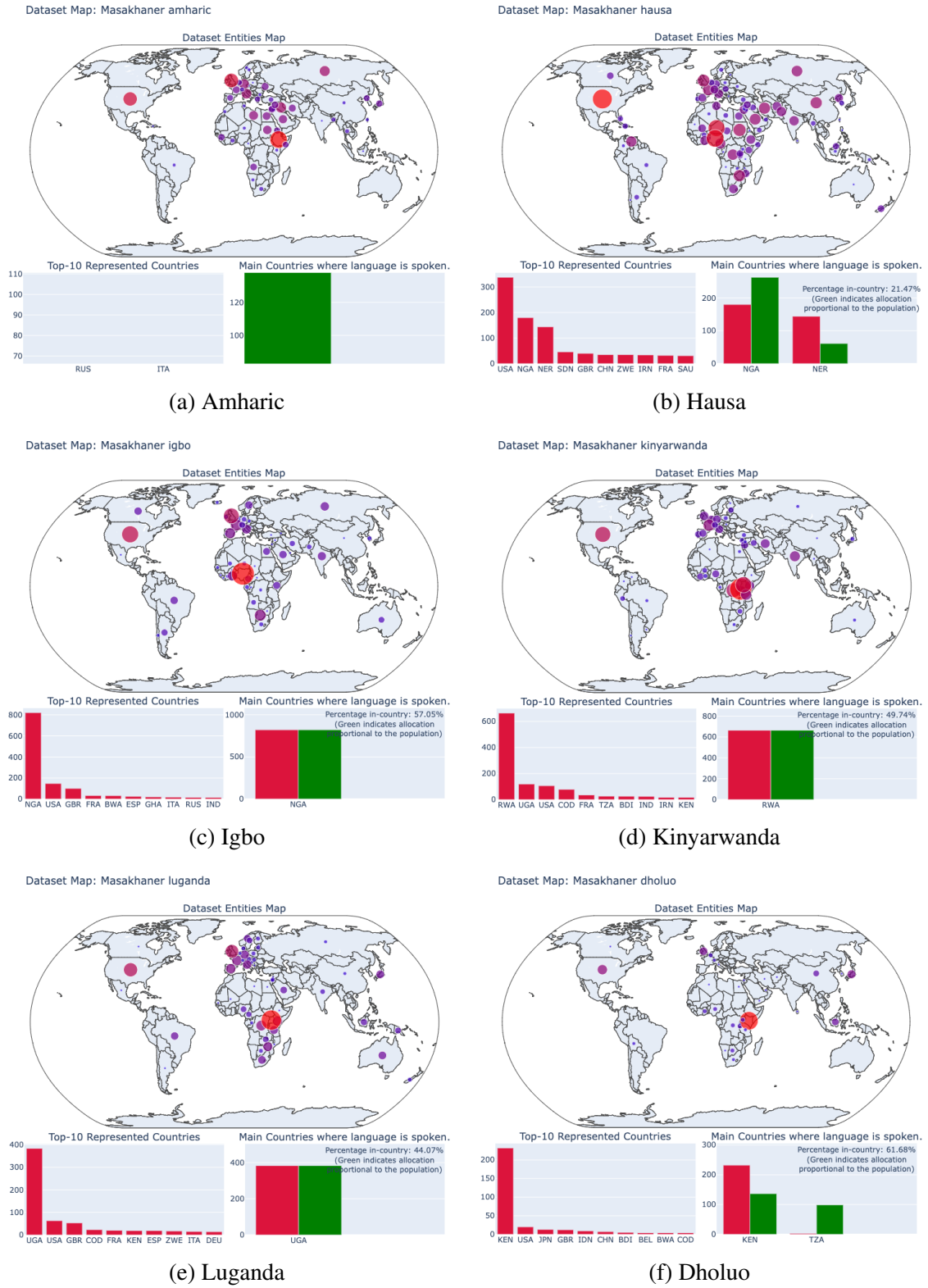
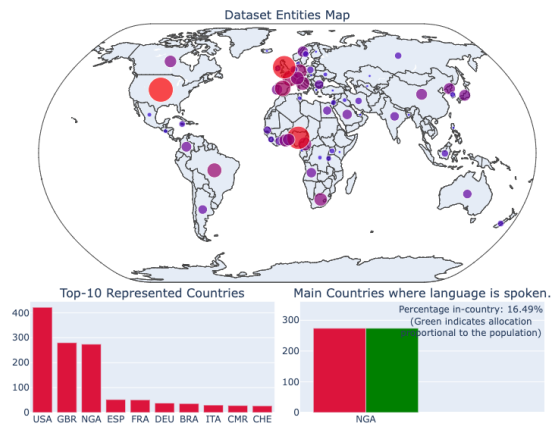


Figure 8: MasakhaNER Geographic Distributions (Part 1).

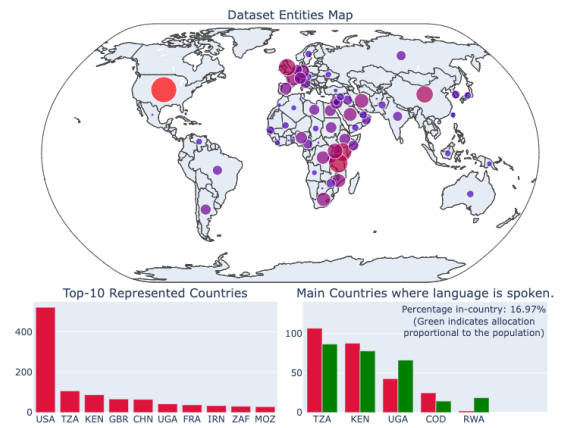
## MasakhaNER Geographic Coverage

Dataset Map: Masakhaner nigerian pidgin



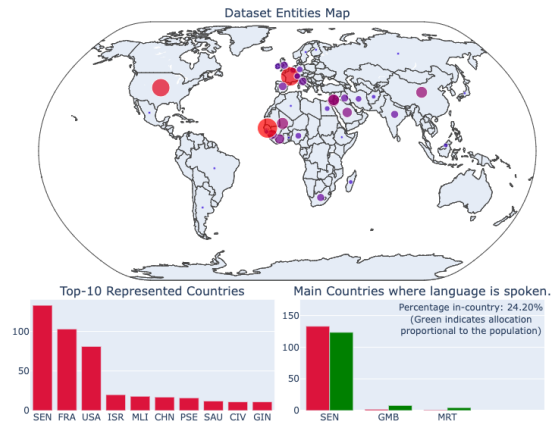
(g) Nigerian English

Dataset Map: Masakhaner swahili



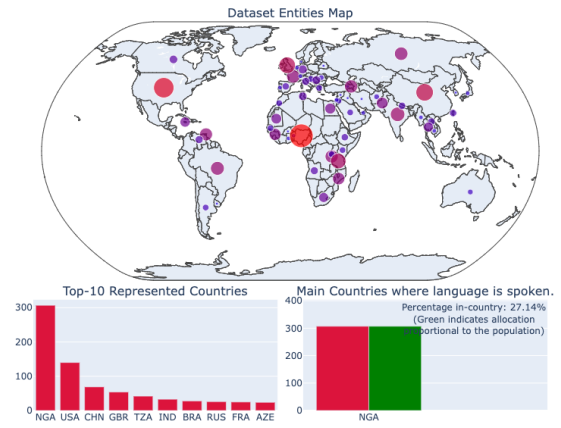
(h) kiSwahili

Dataset Map: Masakhaner wolof



(i) Wolof

Dataset Map: Masakhaner yoruba



(j) Yoruba

Figure 9: MasakhaNER Geographic Distributions (Part 2).

## TyDi-QA Geographic Coverage

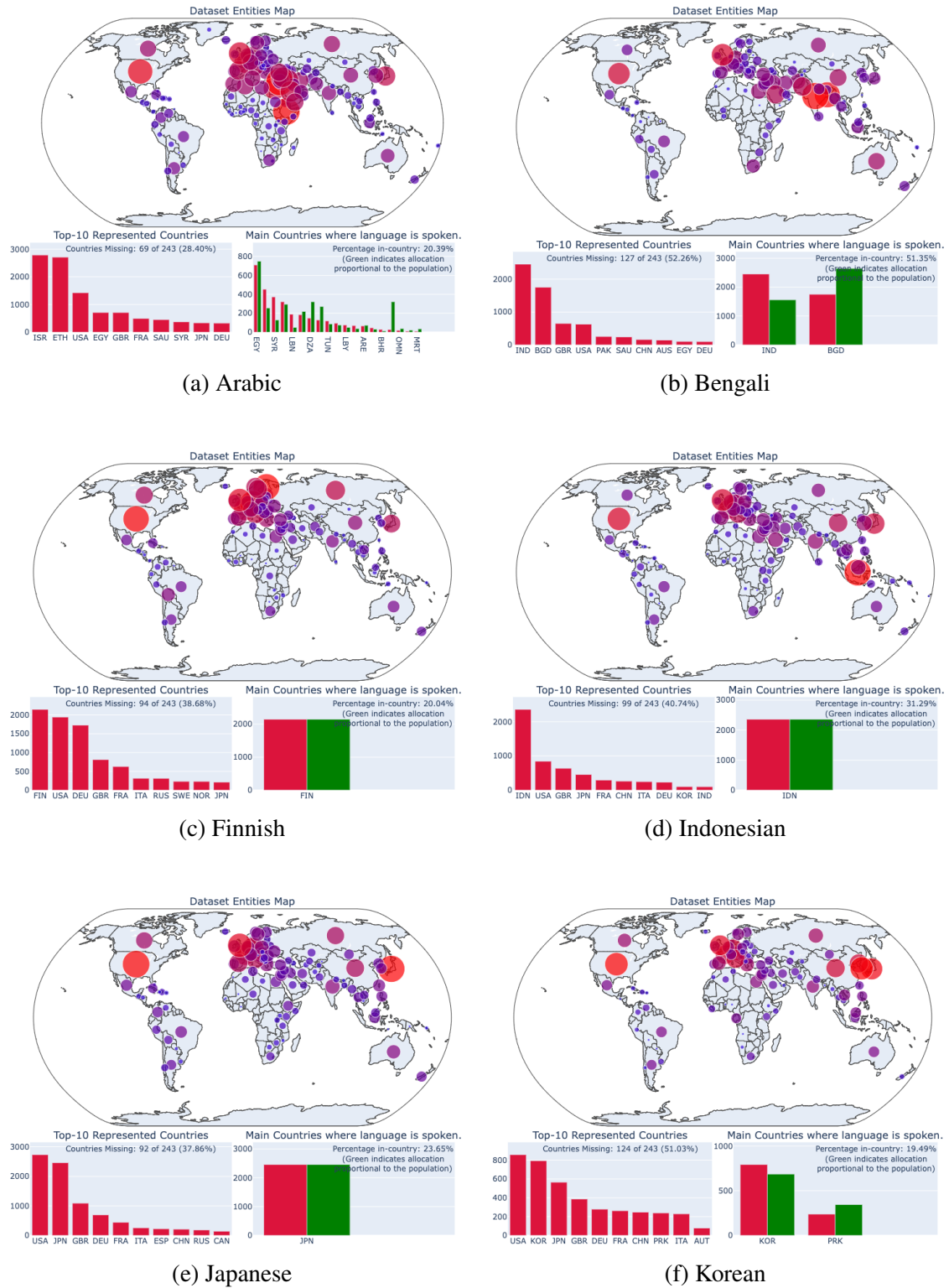


Figure 10: TyDi-QA Geographic Distributions (Part 1).



## TyDi-QA Geographic Coverage

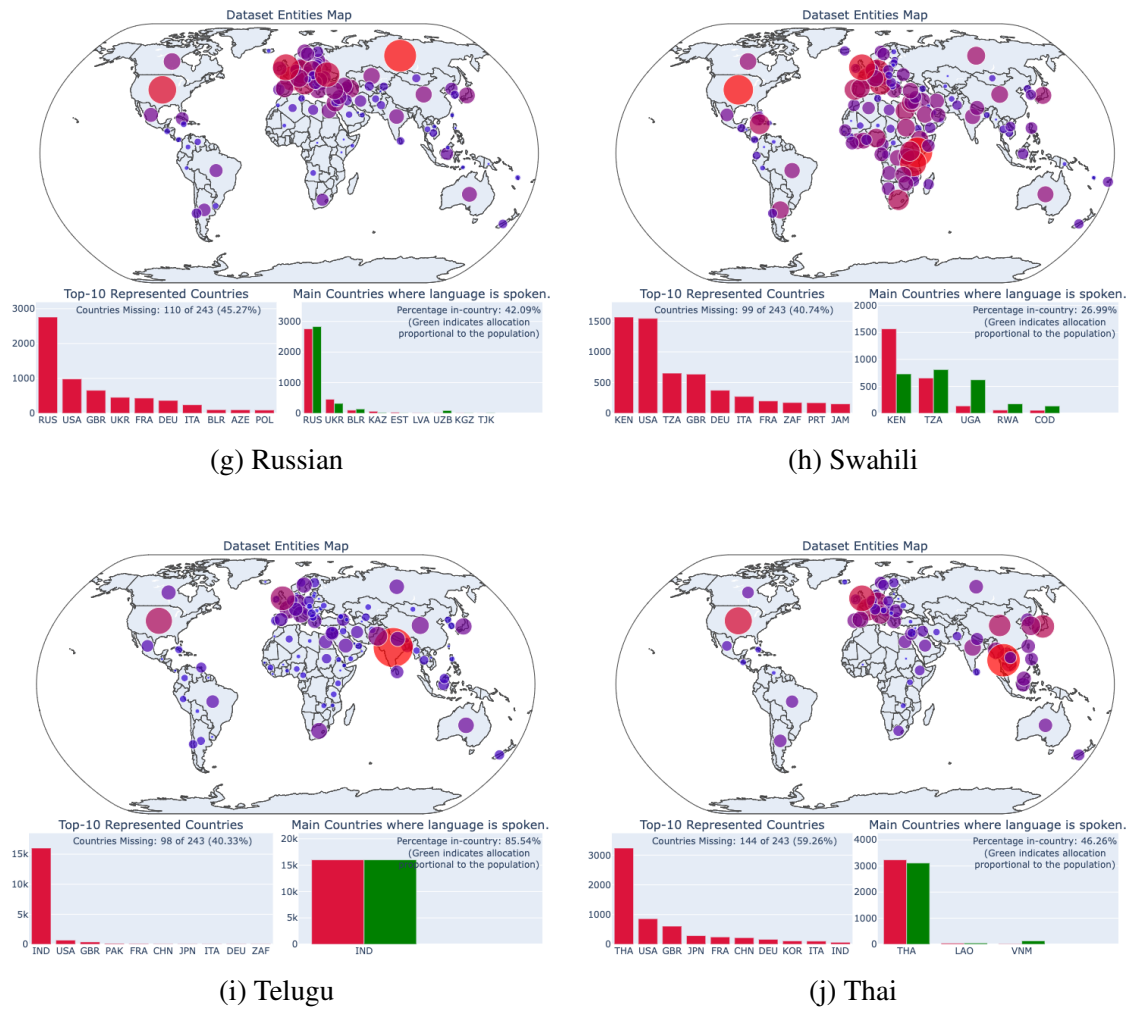


Figure 11: TyDi-QA Geographic Distributions (Part 2).

## Pan-X (WikiANN) Geographic Coverage

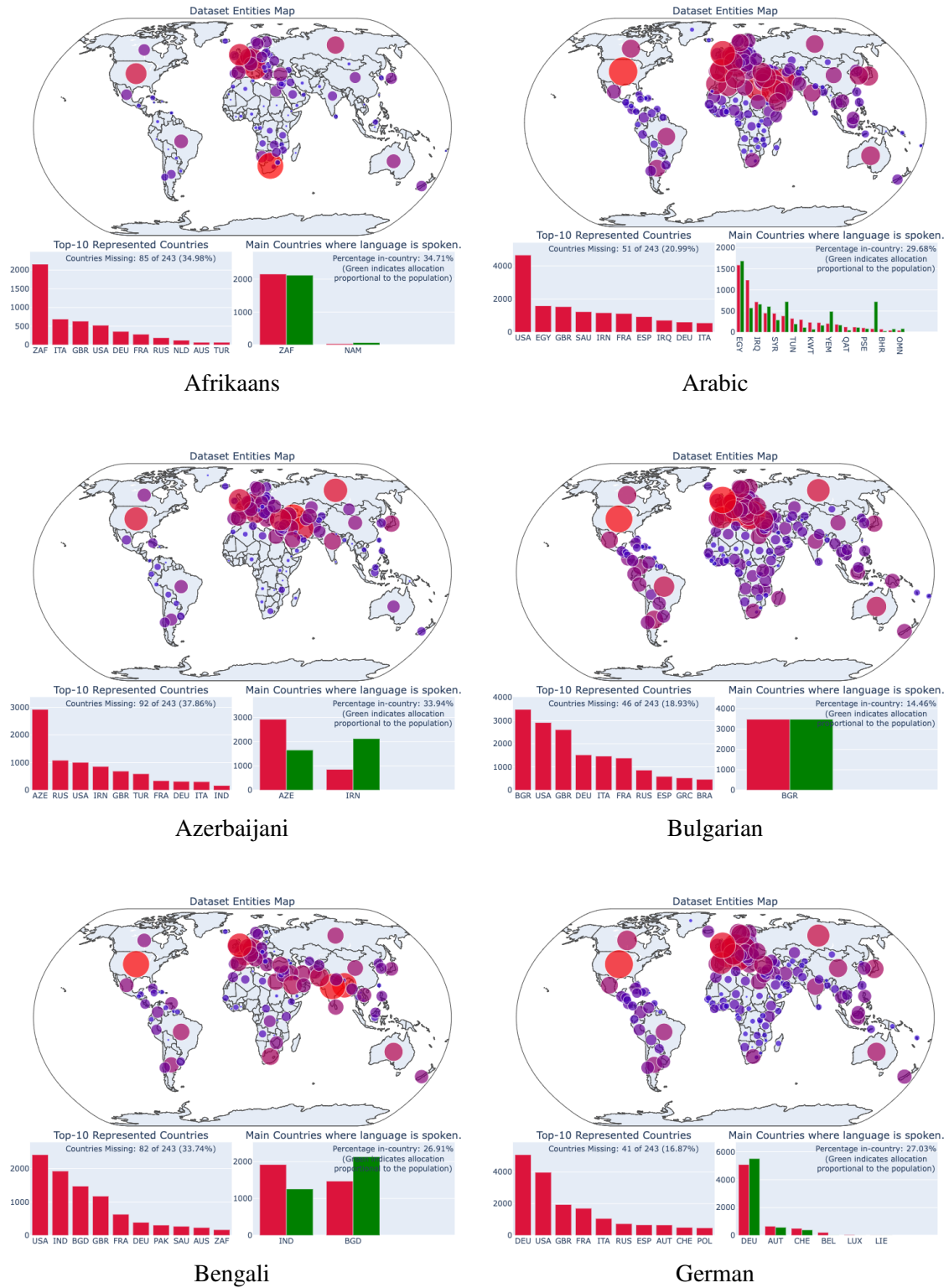


Figure 12: WikiANN Geographic Distributions (Part 1).

## Pan-X (WikiANN) Geographic Coverage

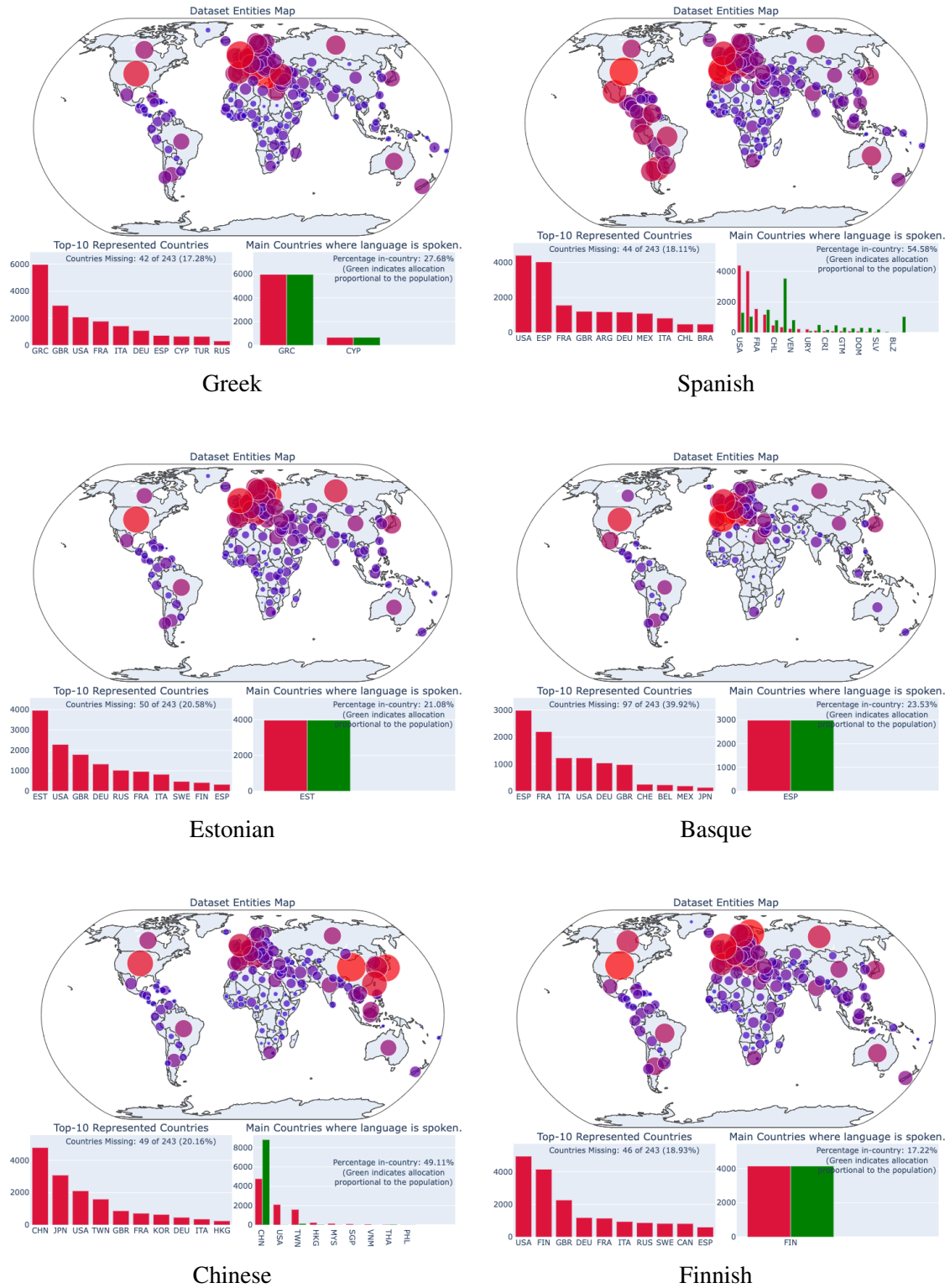


Figure 13: WikiANN Geographic Distributions (Part 2).

## Pan-X (WikiANN) Geographic Coverage

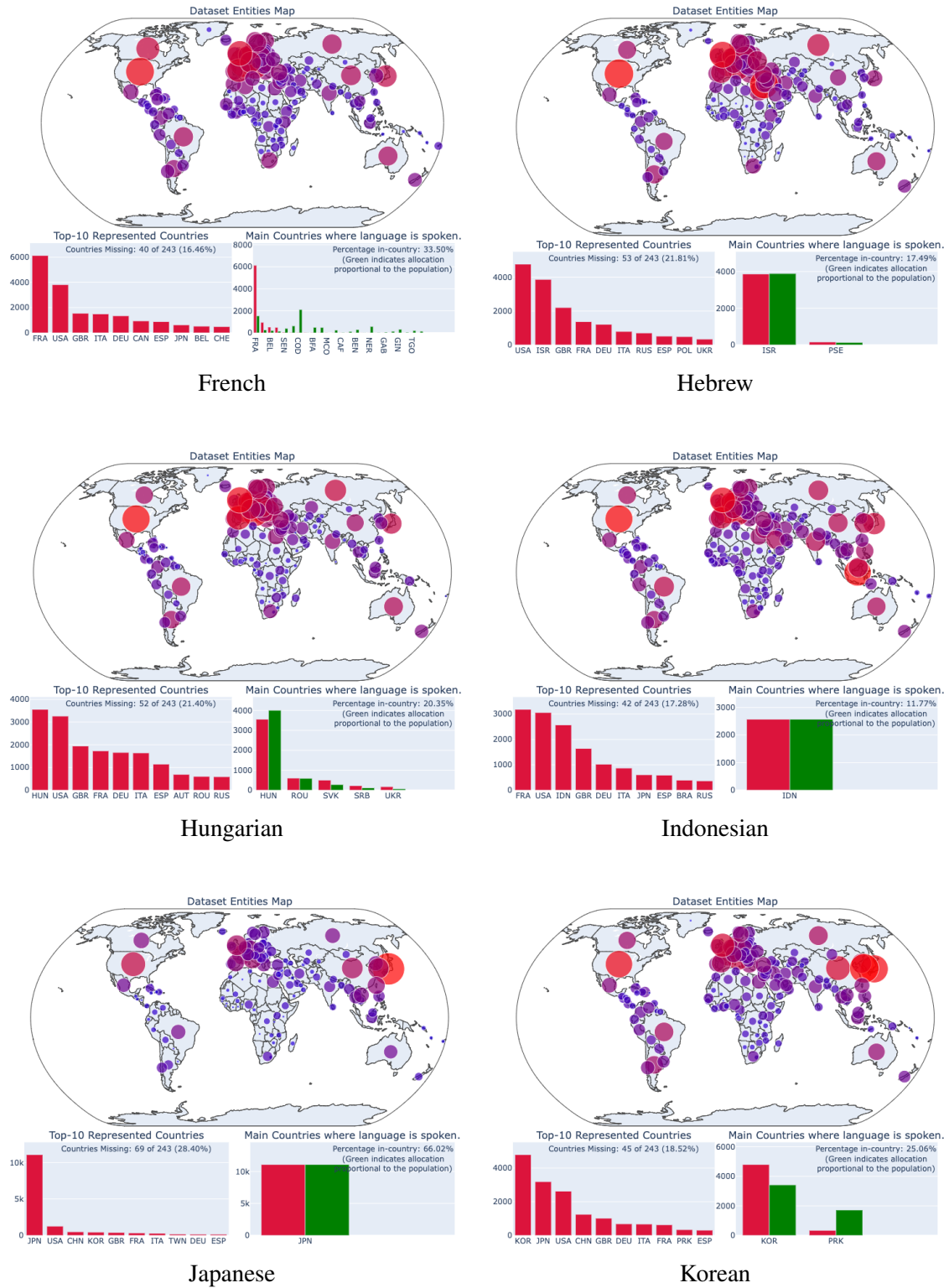


Figure 14: WikiANN Geographic Distributions (Part 3).

## Pan-X (WikiANN) Geographic Coverage

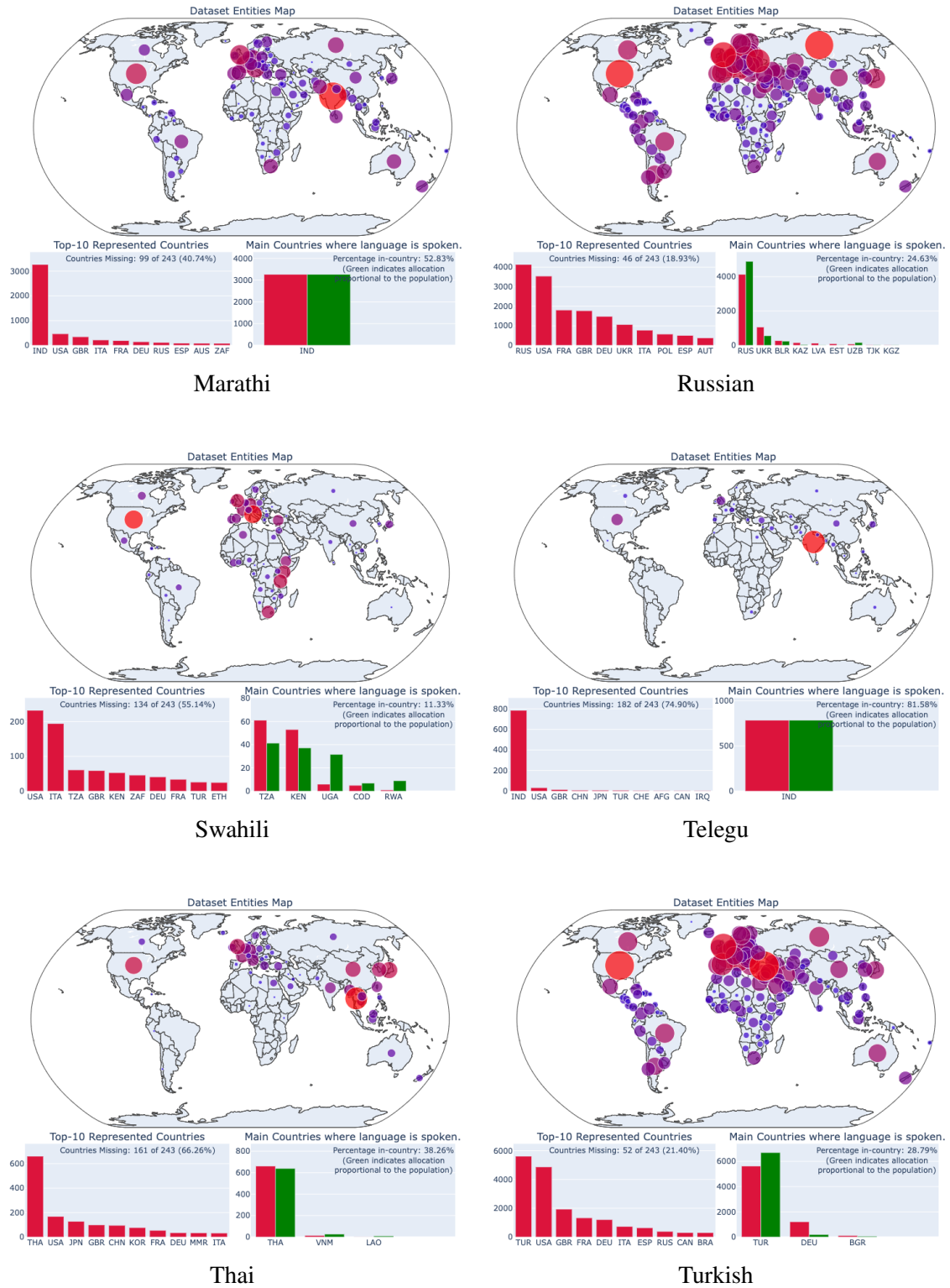
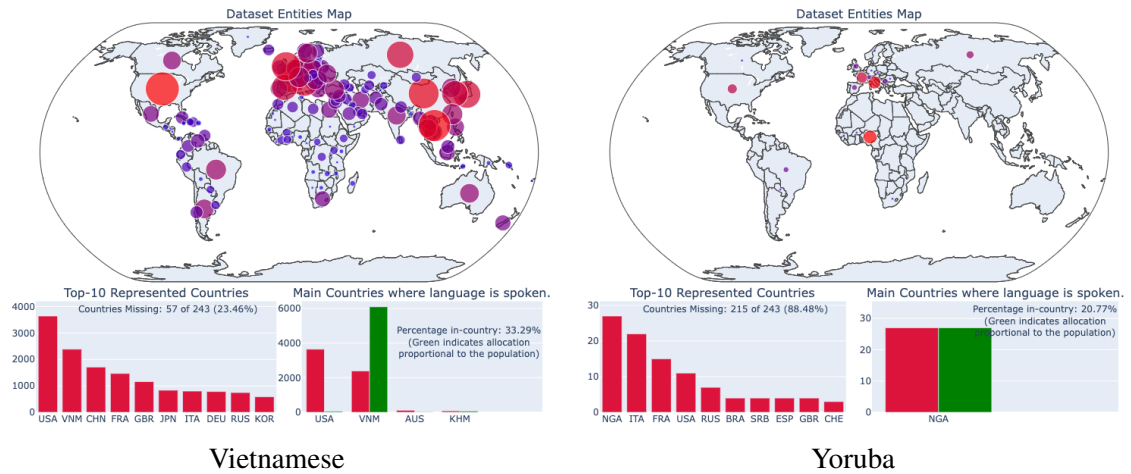


Figure 15: WikiANN Geographic Distributions (Part 4).

## Pan-X (WikiANN) Geographic Coverage



Vietnamese

Yoruba

Figure 16: WikiANN Geographic Distributions (Part 5).

## SQuAD Geographic Coverage

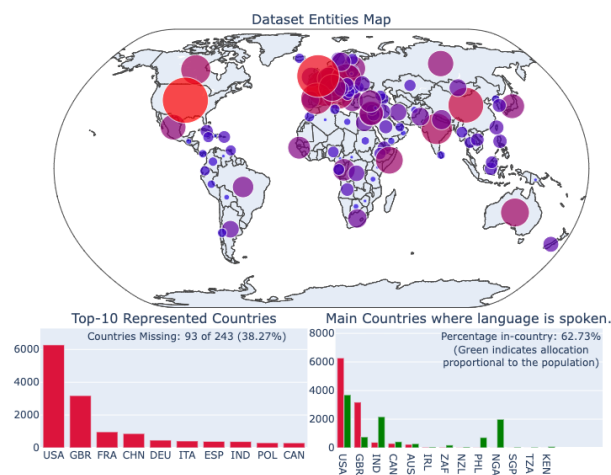


Figure 17: SQuAD Geographic Distributions.

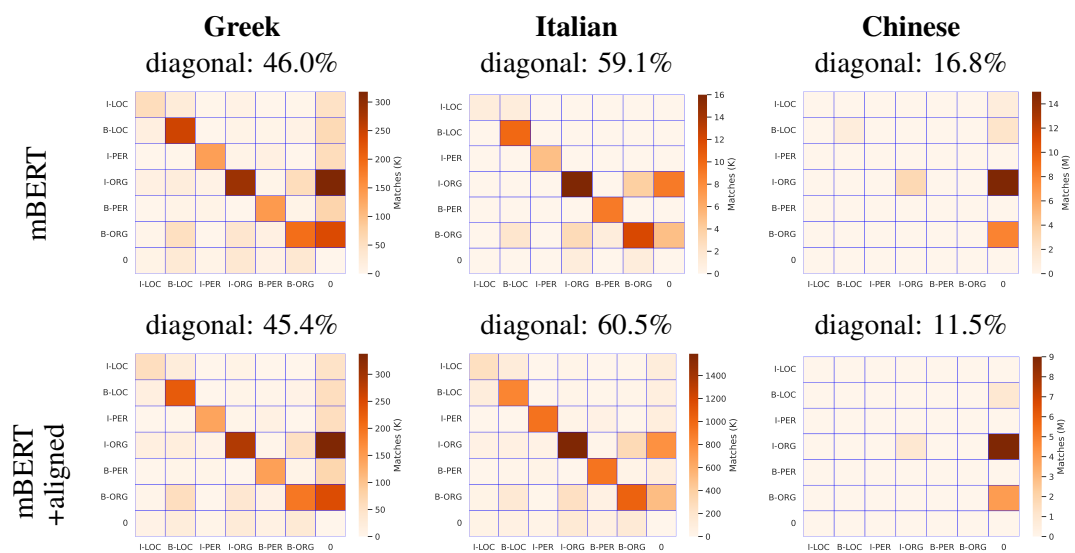


Figure 18: Confusion matrices for Greek, Italian and Chinese.