

---

# EVALUATION OF GEOGRAPHICAL DISTORTIONS IN LANGUAGE MODELS: A CRUCIAL STEP TOWARDS EQUITABLE REPRESENTATIONS

---

A PREPRINT

Rémy Decoupes [remy.decoupes@inrae.fr](mailto:remy.decoupes@inrae.fr)<sup>1,2</sup>, Roberto Interdonato [roberto.interdonato@cirad.fr](mailto:roberto.interdonato@cirad.fr)<sup>1,3</sup>,  
Mathieu Roche [mathieu.roche@cirad.fr](mailto:mathieu.roche@cirad.fr)<sup>1,3</sup>, Maguelonne Teisseire [maguelonne.teisseire@inrae.fr](mailto:maguelonne.teisseire@inrae.fr)<sup>1,2</sup>, and  
Sarah Valentin [sarah.valentin@cirad.fr](mailto:sarah.valentin@cirad.fr)<sup>1,3</sup>

<sup>1</sup>TETIS, Univ. Montpellier, AgroParisTech, CIRAD, CNRS, INRAE. Maison de la Télédétection 500, rue J.F. Breton  
34090 Montpellier  
<sup>2</sup>INRAE  
<sup>3</sup>CIRAD

## ABSTRACT

Language models now constitute essential tools for improving efficiency for many professional tasks such as writing, coding, or learning. For this reason, it is imperative to identify inherent biases. In the field of Natural Language Processing, five sources of bias are well-identified: data, annotation, representation, models, and research design. This study focuses on biases related to geographical knowledge. We explore the connection between geography and language models by highlighting their tendency to misrepresent spatial information, thus leading to distortions in the representation of geographical distances. This study introduces four indicators to assess these distortions, by comparing geographical and semantic distances. Experiments are conducted from these four indicators with ten widely used language models. Results underscore the critical necessity of inspecting and rectifying spatial biases in language models to ensure accurate and equitable representations.

**Keywords** NLP · LLM · Spatial Information

## 1 Introduction

Nowadays Large Language Models (LLMs), generally used through dedicated Chat Bots, have become a primary source of knowledge. Thanks to their effective question-answering system, LLMs are progressively replacing search engines and encyclopedic resources like Wikipedia.

LLMs can be seen as a compression of a large amount of information found on the Internet, and, much of which has a spatial dimension (Manvi et al., 2023). Moreover, questions related to geography, travel, and their cultural aspects represent the third most important use of LLMs, after the ones about programming and artificial intelligence (Zheng et al., 2023). What is more, some kinds of information, such as spatially disaggregated indicators, are not directly accessible through Earth observation and spatial imagery

All this reinforces the use of LLMs for questions related to geography. For instance, the spatial information included in LLMs can also be beneficial for social good-related AI applications, such as crisis management or humanitarian aid (Belliaro et al., 2023).

Since the work by Vaswani et al. (2017), encoder-based language models such as BERT (Devlin et al., 2019), as well as autoregressive models like ChatGPT or LLama2 (OpenAI, 2023; Touvron et al., 2023), have significantly advanced the field of Natural Language Processing (NLP) across a spectrum of tasks, including text classification, named entity recognition, summarization, text generation, and more. These models still undergo the five sources of bias

present in any NLP project (Hovy & Prabhumoye, 2021), namely biases in the data, annotations, vector representations, models and research design. Biases become problematic when they lead to the perpetuation of false ideas. Indeed, inaccurate demographic and cultural representations cause hallucination (Mialon et al., 2023; Puchert et al., 2023) and amplification of well-represented entities at the expense of less-represented ones. These models perpetuate biases originating from their training corpora (Navigli et al., 2023).

In this study, we propose to focus on the biases inherent to spatial information representation. While social biases such as gender, sexual orientation, and ethnicity, are often studied (Navigli et al., 2023), they rarely include biases related to geographical knowledge representation. These distortions in representation lead to a decrease in performance on downstream tasks, as we have demonstrated in our prior research (Decoupes et al., 2023). Indeed, in a text classification task within the context of disaster reporting, fine-tuned models tended to overemphasize associations between locations and type of extreme events, such as "Pakistan" and "Flood". The ability of the model to generalize (e.g., when faced with similar disasters occurring in different locations) can thus be limited by an overfitting phenomenon. In this case, a convenient way to reduce overfitting is to resort to data augmentation (Bayer et al., 2022), i.e., by replacing locations in the original sentences. Notably, we found that replacing distant locations for the original mentions (e.g. replacing "Islamabad" by "Buenos Aeres") significantly improved performance, while the replacement with proximate locations (e.g. replacing "Islamabad" by "New Delhi") yielded marginal improvements. However, the remaining question is: why is a distant location preferable to a closer one?

In this paper, we propose to identify these misrepresentations to pinpoint both isolated and over-represented areas. To assess the reliability of geographical knowledge in different language models, we introduce four criteria: basic geographical knowledge assessment, geographical information significance in training datasets, geographical and semantic distances level of correlation, and anomalies observed. These criteria are then employed to evaluate ten of the most commonly utilized models.

## 2 Related work

In recent years, NLP has been revolutionized by the arrival of attention-based models and their transformer layers (Vaswani et al., 2017). Nowadays, a large number of pre-trained language models have been deployed, that are either generalist or specialized to a given task and application domain.

To build their representations of a language, i.e. the link between words and their numerical vectors (embeddings), language models are trained in a self-supervised mode on huge corpora of texts. BERT was trained on 2.5 billion words from Wikipedia and 0.8 billion words from Google Books (Devlin et al., 2019), while Llama2 has been trained mainly on Common-Crawl, a web dataset containing 1,000 billion words (Touvron et al., 2023). The first difference between the encoder-based models and the generative models (LLMs) is the task used for their pre-training. BERT uses self-masking: the model masks itself random words and tries to predict them to adjust the weights of its neurons. LLMs, on the other hand, are trained to predict the next word and then, they are fine-tuned and aligned on multiple datasets.

To have a better impact on downstream tasks, several studies aim to enhance the geographical knowledge of models by injecting spatial information into the questions (or prompts) addressed to them (Hu et al., 2023; Mai et al., 2023). Another way to incorporate geographical knowledge into models is by modifying the neural network architecture with a merging layer between the vector representations (embeddings) and geographical knowledge graphs, enabling the enhancement of representations (Huang et al., 2022; Li et al., 2023).

However, the reliability of the spatial knowledge in LLMs is not uniform, for instance, GPT-4 shows imprecision when asked to provide GPS coordinates or distances between sparsely populated cities (Roberts et al., 2023). In particular, the knowledge is not worldwide uniform, some areas seem to be underrepresented, thus directly impacting the performance of these models on geography-related NLP tasks, such as crucial tasks in humanitarian crisis response (Belliaro et al., 2023). This is what we will explore in the next sections.

## 3 Geographical knowledge reliability indicators

To assess the biases of spatial representation, we propose four indicators with two main objectives. The primary objective is to assess the reliability of geographical knowledge to identify potential regions of the world that are less well-known to the models. The second objective is to assess whether the representations of places are spatially biased by examining the relationships between geographical and semantic distances. As listed in Tab. 1, we include two language model families, i.e., encoder-based models (e.g. BERT) and LLMs (e.g. ChatGPT). For encoder-based models, we include the most known models and their multilingual versions: BERT (Devlin et al., 2019), BERT-multilingual, RoBERTa (Liu et al., 2019), XLM-RoBERTa (Conneau et al., 2020) and a BERT-based geospatial understanding model,

Models	Name	Languages	Parameters (millions)	Vocabulary Size
bert	bert-base-uncased	English	110	30522
bert multi-lingual	bert-base-multilingual-uncased	6	168	105879
geolm	geolm-base-cased	NA	28996	NA
roberta	roberta-base	English	125	50265
xlm-roberta	xlm-roberta-base	100	279	250002
llama2	Llama-2-7b-chat-hf	English	7000	32000
mistral	Mistral-7B-Instruct-v0.1	English	7000	32000
ada	text-embedding-ada-002	NA	NA	NA
chatGPT	gpt-3.5-turbo-0301	NA	NA	NA

Table 1: Parameters of the Language Models included in the study. NA (Not Available) is used when the parameters could not be found in the associated paper.

GeoLM (Li et al., 2023). Concerning LLMs, we compare four of them by including two among the most reused open source models<sup>1</sup>: Llama-2 (Touvron et al., 2023) & Mistral (Jiang et al., 2023) and OpenAI/ChatGPT will be utilized with prompts, while OpenAI/Ada will be employed for retrieving embeddings due to the unavailability of ChatGPT embedding. As languages significantly contribute to cultures and consequently to geographical representations, we have included multilingual models like XLM-RoBERTa and BERT-multilingual in our comparisons. As 90% of Llama-2 and Mistral training data are in English (followed by programming languages), we consider them to be English language models. Furthermore, even though ChatGPT can respond in multiple languages, we are not aware of the language distribution during its training. Therefore, we cannot assert whether it is a multilingual model or predominantly English. For ground truth, we use geographical data sourced from the GeoNames gazetteer extracted by OpenDatasoft, comprising cities with a population of over 1000 inhabitants<sup>2</sup>.

### 3.1 Disparity of geographical knowledge reliability indicators

Two indicators are proposed to assess the disparity of geographical knowledge in language models across all regions of the world. By geographical knowledge, we refer to the model’s ability to provide knowledge about human geography, such as political boundaries.

#### 3.1.1 Indicator 1: spatial disparities in geographical knowledge

This indicator aims to assess geographical knowledge by evaluating the models’ ability to predict the country based on its capital. To calculate it, we had to adapt the probe to the two families of models (encoder-based and LLMs). Encoder-based models were initially pre-trained to predict a masked word in a sentence. Based on the context of the masked word, these models can guess it (illustrated by listing 1). Since LLMs are generative models, we simply propose to ask them the question in natural language in a prompt (illustrated by listing 2).

```
1 masked_sentence = f'{{city}} is capital of <mask>'.
```

Listing 1: predicting masked country for encoder-based models

```
1 {"role": "user", "content": "Name the country corresponding to its capital: Paris.
   Only give the country."},
2 {"role": "assistant", "content": "France"},
3 {"role": "user", "content": f"Name the country corresponding to its capital: {{city}}.
   Only give the country."}
```

Listing 2: question answering for LLMs

<sup>1</sup>[https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard)

<sup>2</sup><https://public.opendatasoft.com/explore/dataset/geonames-all-cities-with-a-population-1000>

### 3.1.2 Indicator 2: Spatial information coverage in training datasets

The main source of bias comes from the quality of the training datasets (Navigli et al., 2023). Therefore, for this indicator, we propose to indirectly assess the spatial coverage of the training datasets. To do this, we examine the vocabulary of the models that correspond to the most frequent words encountered during their pre-training. We therefore look at the number of cities by continent that appear in the list of words most frequently encountered. Consequently, this experiment seeks to identify city names present in the vocabularies of language models, indirectly gauging the over- or under-representation of these cities by continent in the training datasets. We are investigating the inclusion or exclusion of cities (using their English names) with a population of over 100,000, totalling 4,916 cities, in the predefined vocabularies of the models. For example, London or Paris are in the BERT vocabulary, but Ouagadougou (capital of Burkina Faso) is not.

## 3.2 Geographical distance distortion indicators

In this section, we present two indicators that aim to assess the correlation between geographic distances and semantic distances for pairs of locations. The models convert the words they encounter into high-dimensional vectors (768 for BERT, 4096 for Llama2). These representations, also known as embeddings, help capture the semantics of words. The embedding of a word is partly formed during the training phase by co-occurrence with other words encountered in the training set, and also during the inference phase, during which the pre-trained embedding is modified by the context of its sentence. According to Gurnee & Tegmark (2023), these embeddings also contain spatial information. Indeed, by training a small model (multi-layer perceptron), they manage to correctly predict the GPS coordinates of locations from their embeddings. We rely on these results to introduce the next two indicators.

### 3.2.1 Indicator 3: Correlation between geographic distance and semantic distance

This indicator aims to analyze whether semantic representations (embeddings) are correlated with geographical distances. The objective is to assess, continent by continent, whether semantic representations take into account the geographical distance between pairs of cities. Geographical distance represents the direct distance between two cities as the crow flies. Among all semantic distances used in NLP (Azarpanah & Farhadloo, 2021), we opted for the one based on cosine similarity because it is the most widely used. Semantic distance  $D_{sem}$  is defined as the complement of the cosine similarity between the vectors corresponding to the cities in the embedding space:

$$D_{sem} = 1 - \text{cosine\_similarity}(\text{Embedding}_{\text{city}_1}, \text{Embedding}_{\text{city}_2}) \quad (1)$$

### 3.2.2 Indicator 4: Anomaly between geographical distance and semantic distance

This indicator aims to identify regions that are semantically distant from the rest of the world in terms of vector average, which could be attributed to an under-representation of these regions in the training data. To visualize these differences, we propose focusing only on the top 3 cities per country and calculating the average semantic distances separating them from all the worldwide other cities under consideration.

We therefore apply this indicator to all countries. To better measure the semantic isolation, we also correct it for geographical isolation. For example, the average semantic distances from Sydney, Australia to other cities in the world need to be corrected to compare them with the average semantics of a city in Europe since Australia is a geographically distant country. That is why, we propose the metric GDI (Geographical Distortion Index), which is defined as follows

$$GDI = \frac{1 + D_{sem}}{1 + \overline{D_{geo}}} \quad (2)$$

with  $D_{sem}$ : Semantic distance between pair of cities,  $\overline{D_{geo}}$ : Normalized geographic distance among all the cities  $\in [0, 1]$ .

## A four indicator-based comparison

We propose to apply these four indicators as a basis for comparison to quantify the geographical biases of the models we are evaluating. The results are presented and discussed in the next section.

## 4 Results

This section details the results obtained from the four indicators previously mentioned for the ten widely used models outlined in Tab. 1.

	N.Am	S.Am	Eur.	Africa	Asia	Ocea.	World
bert	<b>35.14</b>	78.57	74.0	55.56	<b>75.47</b>	<b>16.67</b>	<b>57.94</b>
bert multilingual	32.43	<b>85.71</b>	<b>78.0</b>	<b>61.11</b>	66.04	<b>16.67</b>	<b>57.94</b>
geolm	0.00	0.00	0.0	0.00	0.00	0.00	0.00
roberta	16.22	57.14	60.0	18.52	54.72	4.17	36.05
xlm-roberta	13.51	50.00	50.0	7.41	35.85	0.00	25.75
llama2	64.86	<b>92.86</b>	<b>94.0</b>	83.33	<b>84.91</b>	62.50	81.12
mistral	40.54	57.14	54.0	55.56	60.38	29.17	51.07
chatgpt	<b>75.68</b>	85.71	82.0	<b>90.74</b>	83.02	<b>75.00</b>	<b>82.40</b>
Nb countries	37	14	50	54	53	24	232

Table 2: Percentage of correct predictions of the country given its capital, averaged among continents. The abbreviations correspond as follows: *N. Am*: North America, *S. Am*: South America, *Eur.*: Europe, *Ocea.*: Oceania. For each continent, the best results by model family (encoder-based and LLMs) are in bold.

#### 4.1 Disparity of geographical knowledge representations

Here we present the results of the first two experiments: evaluation of basic geographical knowledge (link between countries and their capital) and most encountered city names in the training dataset.

##### 4.1.1 Indicator 1: spatial disparities in geographical knowledge

To evaluate the geographical knowledge, models are compared through the task of predicting the country from its capital. Fig. 1 illustrates the correct answers for BERT and ChatGPT in the form of a map. Tab. 2, on the other hand, provides results in terms of percentages for all models.

The first observation is that the reliability of geographical knowledge is not directly related to the number of parameters in the models. For instance, BERT (110 million) achieves more accurate predictions than Mistral (7 billion). The training corpora seem to play a more significant role in this capability. BERT, which was trained on 2.5 billion words from Wikipedia and 800 million from Google Books, was, we suppose, trained with this type of geographical knowledge. Mistral may be less exposed to geographical information during its training. Although its architecture gives it greater information retention capacity, it remains less effective than BERT (with a difference of -8 points). This suggests that the bias introduced by the training dataset is one of the biases that have the most impact on downstream performance (Navigli et al., 2023).

While the poor results in North America and Oceania can be related to the fact that both have fewer countries (so a poor prediction leads to a significant penalty), Africa appears to be less familiar to the models excepted for ChatGPT, as illustrated in Fig. 1a.

Another interesting observation is that multilingual models, which have been able to capture a greater diversity of cultures, do not necessarily improve the base models. XLM-RoBERTa even performs worse than RoBERTa. However, Bert multilingual provides greater granularity for the European and African continents at the expense of Asia.

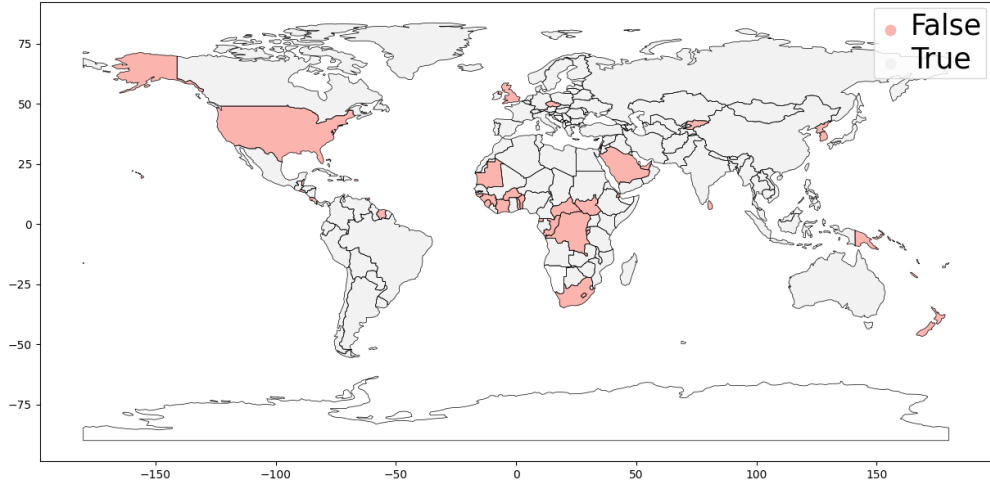
Finally, the results of GeoLM are at 0. We contacted the authors<sup>3</sup> but we did not receive a response.

##### 4.1.2 Indicator 2: Spatial information coverage in training datasets

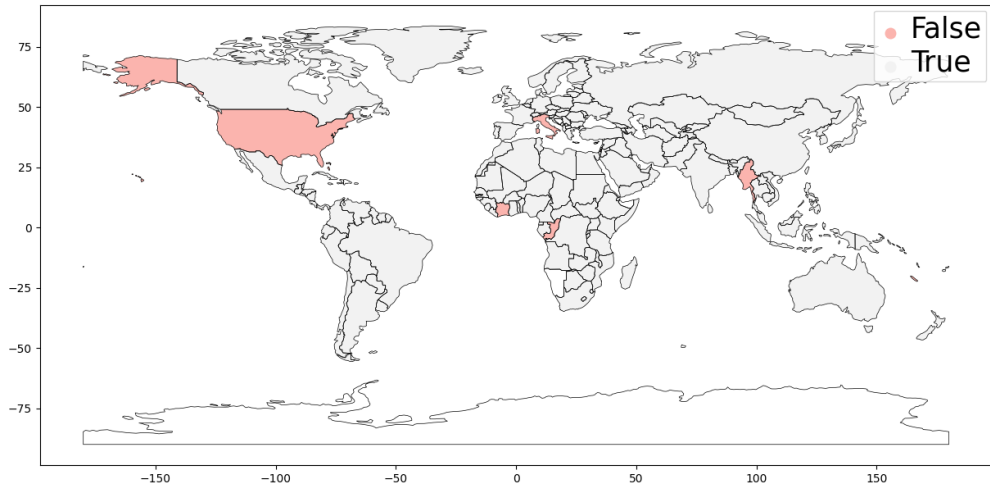
To validate the observations derived from the previous indicator, which suggest that the quality of the training set has a greater impact than the size of the model, we propose to indirectly assess the type of training corpus by counting the number of cities with over 100,000 inhabitants present in the fixed vocabulary of the models. These vocabularies are constructed from the most frequently encountered words during training. To do this, we look at how many cities per country with more than 100,000 inhabitants are in the list of the most frequently encountered words by the models. As with the first experiment, the results for all models are provided in Table 3 and a map illustrates these results for BERT and BERT multilingual in Figure 2. Note that we have used the OpenAI Ada model instead of ChatGPT (GPT-3.5) because its embeddings are not accessible.

The first interesting observation concerns the encoder-based models. BERT and BERT Multilingual, which had better geographical knowledge (with indicator 1), do indeed have the most cities in their vocabularies. Another interesting result is that Bert Multilingual’s training allowed it to find more cities outside the English-speaking world, which on the

<sup>3</sup><https://huggingface.co/zekun-li/geolm-base-cased/discussions/2>



(a) BERT



(b) ChatGPT (GPT-3.5-turbo-0301)

Figure 1: Correct prediction of the country given its capital

other hand led to a loss of coverage for Oceania and North America. This is also in line with the results of indicator 1. Bert Multilingual performs better on non-English speaking continents thanks to its more diverse training corpus.

However, when it comes to LLMs, it’s impossible to analyze the spatial coverage of the training data sets by working only at the level of their vocabularies. There are very few cities in their vocabularies. Moreover, according to manual analysis, the cities that appear in the most frequent words of the models are cities that are spelt like common words, such as Bath (United Kingdom).

## 4.2 Geographical distance distortion

In this section, we present the correlation between geographical and semantic distances for all regions of the world. During their training, models learn word representations based on other words that co-occur in sentences. The following two indicators aim to analyze whether this training method enables the models to represent geographical distances between spatial entities.

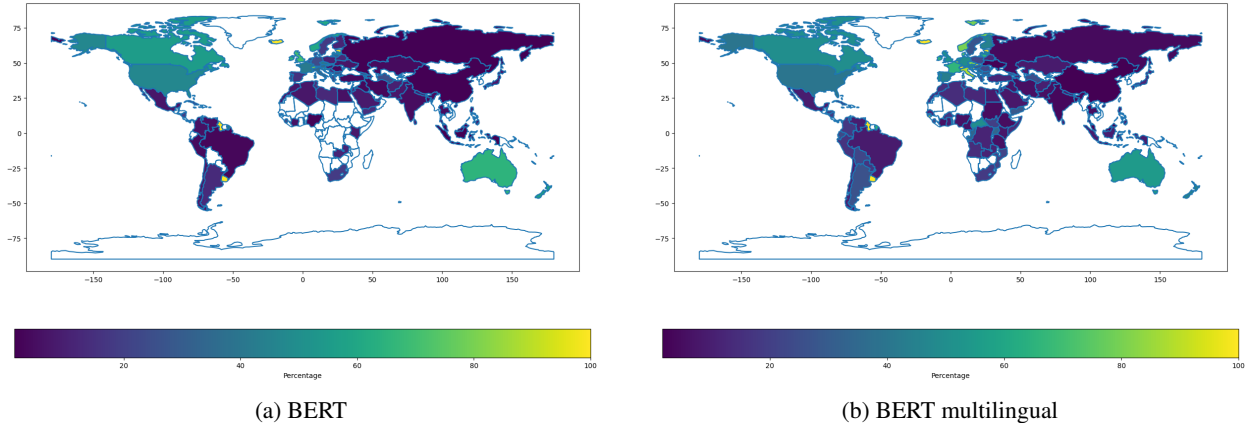


Figure 2: Percentage of cities of more than 100k inhabitants in models vocabulary. Countries in white have none of their cities in the vocabulary. The best results by model family (encoder-based and LLMs) are in bold.

	N.Am	S.Am	Eur.	Africa	Asia	Oceania
bert	<b>33.80</b>	4.82	22.88	4.98	4.50	<b>60.0</b>
bert mlingual	30.05	<b>13.00</b>	<b>36.05</b>	<b>10.82</b>	<b>6.16</b>	50.0
geolm	1.41	0.00	0.67	0.69	0.26	0.0
roberta	0.63	0.84	1.56	1.37	0.66	0.0
xlm-roberta	1.10	1.47	2.23	3.44	2.62	0.0
llama2	0.31	0.00	0.22	<b>0.34</b>	0.00	0.0
mistral	<b>0.47</b>	0.00	<b>0.33</b>	<b>0.34</b>	0.00	0.0
Nb of cities	639	477	896	582	2289	30

Table 3: Percentage of cities of more than 100k inhabitants in models vocabulary, averaged by continent. The abbreviations correspond as follows: *N. Am.*: North America, *S. Am.*: South America, *Eur.*: Europe. The best results by model family (encoder-based and LLMs) are in bold.

#### 4.2.1 Indicator 3: Correlation between geographic distance and semantic distance

We compare the correlation between semantic and geographical distances pairwise for cities with a population of over 1 million. Figure 3 illustrates this type of result for the BERT model for Europe. We apply a linear regression between geographic and semantic distance and display the  $R^2$  of the linear regression for all models in Table 4.

The first observation is that the correlation coefficients ( $\in [0, 1]$ ) are low for all models. This indicates that the models’ representations (or embeddings) struggle to accurately capture the geographical distances between locations, as suggested by (Decoupes et al., 2023). However, for the encoder-based models, GeoLM is competitive with the best LLM in this task. GeoLM’s specific training from sentences generated from OpenStreetMap, containing neighbours of locations, allows it to grasp geographical distances better. Once again, the training dataset is crucial.

Regarding LLMs, Llama2 did not capture geographical distances, whereas OpenAI/ada, OpenAI’s embedding model, achieves the highest scores. The final observation is that Europe, Asia, and to a lesser extent North America are the continents whose geographical distances are best captured. However, these continents, in the arrangement of their countries, do not share the same characteristics. Indeed, Europe is the smallest continent with the most small countries, while North America is a continent that contains large countries.

Thus, the characteristics of continents do not seem to explain the observed differences between them. Two factors can influence a low correlation. The first is that the location is overrepresented in the training data, which leads to a smaller semantic distance to other locations. The second is the opposite: the location is underrepresented in the training data, so its embeddings are more distant from the embeddings of other locations. Indicator 4 determines whether countries are distant or in the center of the semantic space.

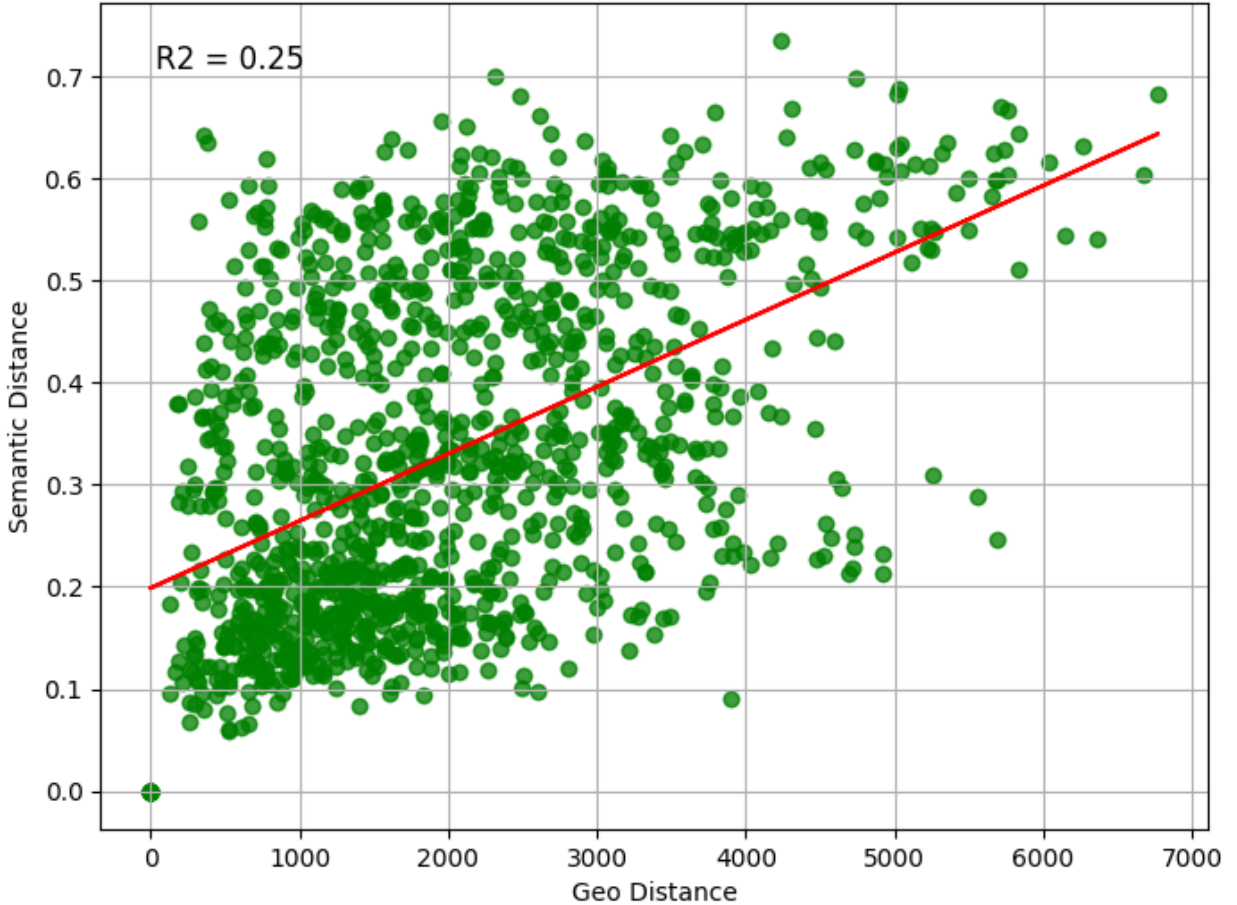


Figure 3: Linear regression between geographical distance (km) and semantic distance with BERT for Europe

	N.Am.	S.Am.	Europe	Africa	Asia	Oceania
bert	0.08	<b>0.07</b>	0.25	0.07	0.09	0.03
bert mlingual	0.01	0.01	0.14	0.03	0.02	0.00
geolm	<b>0.17</b>	0.03	<b>0.37</b>	<b>0.09</b>	<b>0.32</b>	<b>0.07</b>
roberta	0.03	0.02	0.13	0.02	0.05	0.00
xlm-roberta	0.05	0.00	0.19	0.01	0.00	0.00
llama2	0.04	0.03	0.07	0.06	0.03	0.02
mistral	0.10	0.06	0.20	0.07	0.12	0.04
openai/ada	<b>0.17</b>	<b>0.20</b>	<b>0.28</b>	<b>0.17</b>	<b>0.38</b>	<b>0.17</b>

Table 4:  $R^2$  of the linear regression between geographical distance and semantic distance between pairs of cities. The abbreviations correspond as follows: *N. Am*: North America, *S. Am*: South America. The best results by model family (encoder-based and LLMs) are in bold.

#### 4.2.2 Indicator 4: Anomaly between geographical distance and semantic distance

As shown in indicator 3, some continents, such as Africa and Oceania, have very low correlations between semantic and geographical distances. With this new and final indicator, we propose to explain these low correlations either by an overrepresentation of countries in the training datasets, positioning the places of these countries in the centre of the semantic space, or, on the contrary, by an underrepresentation of these places, isolating them in this space.

Figure 4 shows, for the BERT model, the country average of the semantic distances for its three most populous cities with the most populous cities in the world. In Africa, we find countries that are the most semantically distant, such as



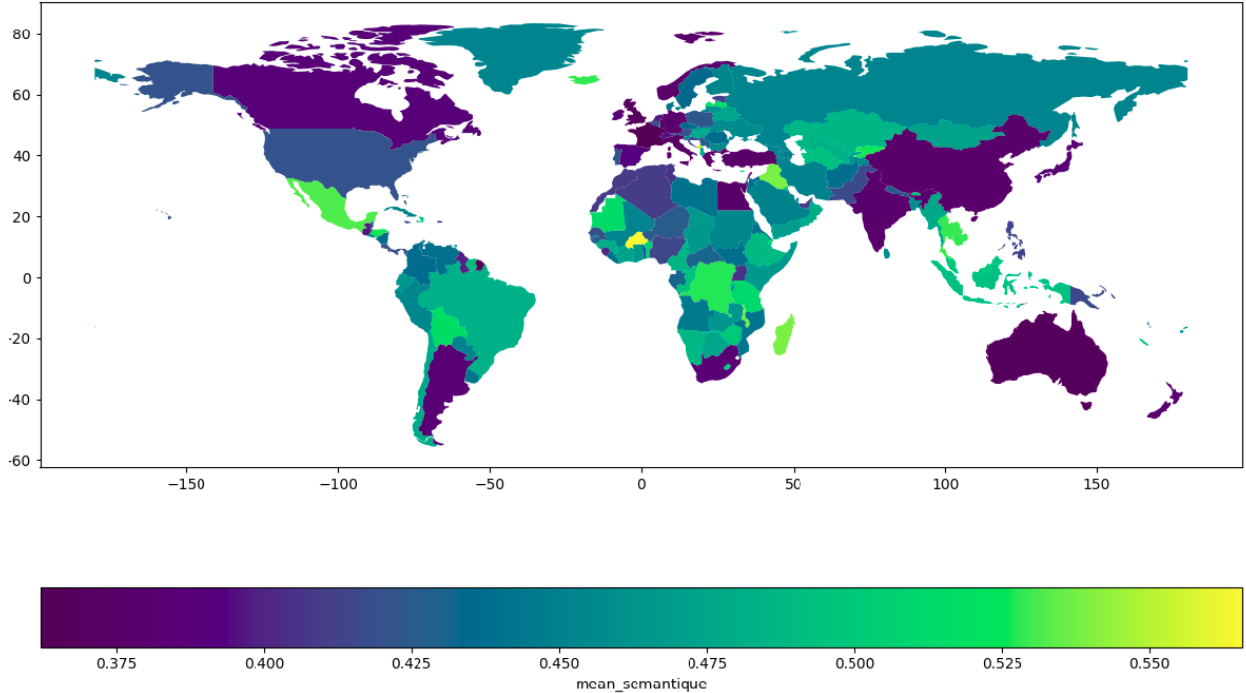


Figure 4: Average per country of the BERT semantic distances from its three most populous cities to the cities of the world

Burkina Faso, the Democratic Republic of Congo or Mauritania. On the other hand, the countries of Oceania, such as Australia and New Zealand, are abnormally close semantically to the most populous cities in the world, compared to their geographical distance.

To validate these observations for all models, we use the GDI ratio (semantic distance divided by geographic distance). Table 5 shows, by model and continent, the average GDI and the number of countries in the 20 most or least distant countries. The first observation is that the trends, described after, are consistent across models. Indeed, for all the models, Europe is the most distant continent to the most populated cities. However as shown by the Figure 4, there are significant disparities between the west and east of Europe, the Eastern European countries are abnormally distant. Then, the African continent appears to be the second most distant. This also confirms the interpretation of the figure, is that Oceania is abnormally close. To a lesser extent, this is also the case for South America.

#### A four indicator-based comparison

To conclude the comparison, four main points should be highlighted. First, a high number of parameters does not guarantee that the model has good geographical knowledge. Indeed, the BERT model (0.10 billion parameters) sometimes outperforms Mistral (7 billion parameters).

The second point is in line with the fact that training datasets are the source of the strongest biases (Navigli et al., 2023). These two observations open up a new perspective of research. It is to analyze the impact of training datasets on the number of parameters and assess the risks of diluting geographical information in favour of better performance in other criteria. This is exemplified by the case of Mistral, which outperforms Llama-7-B and even Llama-13-B in a majority of evaluations (Jiang et al., 2023) but could predict less country given its capital than BERT and Llama-7-B. It is interesting to analyze when geographical information disappears in favour of other knowledge during the training process.

The third point is that multi-language models partially correct the bias of the data by having more diversified training datasets that provide other cultural and geographical contexts, they perform better for African or Asian countries but at the expense of English-speaking countries.

Finally, the last point is that countries in Oceania appear to be abnormally close semantically, while Eastern Europe, Africa and Southeast Asia (Indonesia, Vietnam, Cambodia) are semantically isolated.

Model	Metric	N.America	S.America	Europe	Africa	Asia	Oceania
bert	mean	1	0.99	<b>1.11</b>	1.09	1.07	<u>0.88</u>
	farthest	0	0	<b>11</b>	5	4	0
	nearest	<u>6</u>	4	0	0	4	<u>6</u>
bert-ml	mean	0.87	0.86	<b>0.95</b>	0.94	0.93	<u>0.77</u>
	farthest	1	0	6	6	<b>7</b>	0
	nearest	4	5	0	0	5	<u>6</u>
roberta	mean	0.73	0.72	<b>0.81</b>	0.78	0.77	<u>0.65</u>
	farthest	0	0	<b>15</b>	0	5	0
	nearest	3	<u>7</u>	0	0	4	6
geolm	mean	0.85	0.84	<b>0.96</b>	0.92	0.9	<u>0.75</u>
	farthest	0	0	<b>19</b>	0	1	0
	nearest	5	4	0	0	5	<u>6</u>
roberta-xlm	mean	0.7	0.69	<b>0.78</b>	0.75	0.74	<u>0.62</u>
	farthest	0	0	<b>17</b>	2	1	0
	nearest	3	<u>7</u>	0	0	4	6
llama2	mean	0.96	0.95	<b>1.07</b>	1.04	1.03	<u>0.86</u>
	farthest	0	0	<b>11</b>	4	5	0
	nearest	5	5	0	0	4	<u>6</u>
mistral	mean	0.98	0.96	<b>1.09</b>	1.05	1.04	<u>0.88</u>
	farthest	0	0	<b>15</b>	3	2	0
	nearest	4	<u>6</u>	0	0	4	<u>6</u>
openai/ada	mean	0.84	0.83	<b>0.93</b>	0.9	0.89	<u>0.75</u>
	farthest	0	0	<b>18</b>	1	1	0
	nearest	3	<u>7</u>	0	0	4	6
country		17	12	35	47	38	6

Table 5: GDI ratio (Semantic Distance / Geographic Distance normalized by continent and model). **Mean**: the average GDI per continent. **Farthest**: the number of countries per continent among the 20 farthest (in **bold**). **Nearest**: the number of countries per continent among the 20 least distant (in underline). The abbreviations correspond as follows: *N. America*: North America, *S. America*: South America

This diversity in knowledge could impact NLP tasks related to information retrieval, event detection, and tracking. For example, in crisis management, it is crucial to identify events (nature of the event and location). However, in the massive flow of information exchange, especially through social networks, it is essential to deduplicate events. An event can be described differently, and it is important not to interpret them as several separate events. However, poor semantic representation of locations and distances leads to incorrect deduplication. Therefore, the contribution of these models may have a limited impact in aiding responses to humanitarian crises in countries that may need it the most. This is a second perspective of research that pertains to the correlations between geographical distances (a proxy for cultural context) and semantic distance. How does improving this correlation lead to better performance in downstream tasks such as information retrieval and event detection and tracking? Same question with improving reliability of geographical knowledge, does it leads to better results in geographical NLP tasks for location detection or GPS coordinate assignment ?

## 5 Conclusions

While geography is the third most discussed topic by artificial intelligence and its users, the sources of its biases are underexplored in the literature. In this paper, we propose four indicators to assess the geographical knowledge of language models and their sources of bias. This comparative framework with four indicators is applied to ten models among the most widely used. Significant disparities emerge not only between models but also across continents. Finally, after analyzing these results, we highlight the potential impact on traditional NLP tasks and suggest two avenues for addressing these spatial biases.

## Acknowledgement

This study was partially funded by EU grant 874850 MOOD and is catalogued as MOOD099. The contents of this publication are the sole responsibility of the authors and do not necessarily reflect the views of the European Commission.

## References

- Azarpanah, H., & Farhadloo, M. (2021). Measuring Biases of Word Embeddings: What Similarity Measures and Descriptive Statistics to Use? In *Proceedings of the First Workshop on Trustworthy Natural Language Processing* (pp. 8–14). Online: Association for Computational Linguistics. URL: <https://www.aclweb.org/anthology/2021.trustnlp-1.2>. doi:doi:10.18653/v1/2021.trustnlp-1.2.
- Bayer, M., Kaufhold, M.-A., & Reuter, C. (2022). A Survey on Data Augmentation for Text Classification. *ACM Computing Surveys*, 1, 3544558. URL: <https://dl.acm.org/doi/10.1145/3544558>. doi:doi:10.1145/3544558.
- Belliardo, E., Kalimeri, K., & Mejova, Y. (2023). Leave no Place Behind: Improved Geolocation in Humanitarian Documents. In *Proceedings of the 2023 ACM Conference on Information Technology for Social Good* (pp. 31–39). Lisbon Portugal: ACM. URL: <https://dl.acm.org/doi/10.1145/3582515.3609515>. doi:doi:10.1145/3582515.3609515.
- Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., & Stoyanov, V. (2020). Unsupervised Cross-lingual Representation Learning at Scale. URL: <http://arxiv.org/abs/1911.02116> arXiv:1911.02116 [cs].
- Decoupes, R., Roche, M., & Teisseire, M. (2023). GeoNLPlify: A spatial data augmentation enhancing text classification for crisis monitoring. *Intelligent Data Analysis, Preprint*, 1–25. URL: <https://content.iospress.com/articles/intelligent-data-analysis/ida230040>. doi:doi:10.3233/IDA-230040. Publisher: IOS Press.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North {A}merican Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)* (pp. 4171–4186). Minneapolis, Minnesota: Association for Computational Linguistics volume 1. URL: <https://aclanthology.org/N19-1423>. doi:doi:10.18653/v1/N19-1423.
- Gurnee, W., & Tegmark, M. (2023). Language Models Represent Space and Time. URL: <http://arxiv.org/abs/2310.02207> arXiv:2310.02207 [cs].
- Hovy, D., & Prabhumoye, S. (2021). Five sources of bias in natural language processing. *Language and Linguistics Compass*, 15. URL: <https://onlinelibrary.wiley.com/doi/10.1111/lnc3.12432>. doi:doi:10.1111/lnc3.12432.
- Hu, Y., Mai, G., Cundy, C., Choi, K., Lao, N., Liu, W., Lakhanpal, G., Zhou, R. Z., & Joseph, K. (2023). Geo-knowledge-guided GPT models improve the extraction of location descriptions from disaster-related social media messages. *International Journal of Geographical Information Science*, 37, 2289–2318. URL: <https://www.tandfonline.com/doi/full/10.1080/13658816.2023.2266495>. doi:doi:10.1080/13658816.2023.2266495.
- Huang, J., Wang, H., Sun, Y., Shi, Y., Huang, Z., Zhuo, A., & Feng, S. (2022). ERNIE-GeoL: A Geography-and-Language Pre-trained Model and its Applications in Baidu Maps. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (pp. 3029–3039). URL: <http://arxiv.org/abs/2203.09127>. doi:doi:10.1145/3534678.3539021 arXiv:2203.09127 [cs].
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., & Sayed, W. E. (2023). Mistral 7B. URL: <http://arxiv.org/abs/2310.06825> arXiv:2310.06825 [cs].
- Li, Z., Zhou, W., Chiang, Y.-Y., & Chen, M. (2023). GeoLM: Empowering Language Models for Geospatially Grounded Language Understanding. URL: <http://arxiv.org/abs/2310.14478> arXiv:2310.14478 [cs].
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. URL: <http://arxiv.org/abs/1907.11692> number: arXiv:1907.11692 arXiv:1907.11692 [cs].
- Mai, G., Huang, W., Sun, J., Song, S., Mishra, D., Liu, N., Gao, S., Liu, T., Cong, G., Hu, Y., Cundy, C., Li, Z., Zhu, R., & Lao, N. (2023). On the Opportunities and Challenges of Foundation Models for Geospatial Artificial Intelligence. URL: <http://arxiv.org/abs/2304.06798> arXiv:2304.06798 [cs].
- Manvi, R., Khanna, S., Mai, G., Burke, M., Lobell, D., & Ermon, S. (2023). GeoLLM: Extracting Geospatial Knowledge from Large Language Models. URL: <http://arxiv.org/abs/2310.06213> arXiv:2310.06213 [cs].

- Mialon, G., Dessì, R., Lomeli, M., Nalmpantis, C., Pasunuru, R., Raileanu, R., Rozière, B., Schick, T., Dwivedi-Yu, J., Celikyilmaz, A., Grave, E., LeCun, Y., & Scialom, T. (2023). Augmented Language Models: a Survey. URL: <http://arxiv.org/abs/2302.07842> arXiv:2302.07842 [cs].
- Navigli, R., Conia, S., & Ross, B. (2023). Biases in Large Language Models: Origins, Inventory, and Discussion. *Journal of Data and Information Quality*, 15, 1–21. URL: <https://dl.acm.org/doi/10.1145/3597307>. doi:doi:10.1145/3597307.
- OpenAI (2023). ChatGPT [Large language model]. URL: <https://chat.openai.com>.
- Puchert, P., Poonam, P., van Onzenoodt, C., & Ropinski, T. (2023). LLMMaps – A Visual Metaphor for Stratified Evaluation of Large Language Models. URL: <http://arxiv.org/abs/2304.00457> arXiv:2304.00457 [cs].
- Roberts, J., Lüddecke, T., Das, S., Han, K., & Albanie, S. (2023). GPT4GEO: How a Language Model Sees the World’s Geography. URL: <http://arxiv.org/abs/2306.00020> arXiv:2306.00020 [cs].
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., & al. (2023). Llama 2: Open Foundation and Fine-Tuned Chat Models. URL: <http://arxiv.org/abs/2307.09288>. doi:doi:10.48550/arXiv.2307.09288 arXiv:2307.09288 [cs].
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention Is All You Need. URL: <http://arxiv.org/abs/1706.03762> number: arXiv:1706.03762 arXiv:1706.03762 [cs].
- Zheng, L., Chiang, W.-L., Sheng, Y., Li, T., Zhuang, S., Wu, Z., Zhuang, Y., Li, Z., Lin, Z., Xing, E. P., Gonzalez, J. E., Stoica, I., & Zhang, H. (2023). LMSYS-Chat-1M: A Large-Scale Real-World LLM Conversation Dataset. URL: <http://arxiv.org/abs/2309.11998> arXiv:2309.11998 [cs].