

Multilingual Hallucination Gaps

Cléa Chataigner

Mila, McGill University

CLEA.CHATAIGNER@MILA.QUEBEC

Afaf Taïk

Mila, University of Montréal

AFAF.TAIK@MILA.QUEBEC

Golnoosh Farnadi

Mila, McGill University

FARNADIG@MILA.QUEBEC

Editors: Miriam Rateike, Awa Dieng, Jamelle Watson-Daniels, Ferdinando Fioretto, Golnoosh Farnadi

Abstract

Large language models (LLMs) are increasingly used as alternatives to traditional search engines given their capacity to generate text that resembles human language. However, this shift is concerning, as LLMs often generate hallucinations—misleading or false information that appears highly credible. In this study, we explore the phenomenon of hallucinations across multiple languages in free-form text generation, focusing on what we call *multilingual hallucination gaps*. These gaps reflect differences in the frequency of hallucinated answers depending on the prompt and language used. To quantify such hallucinations, we used the FACTSCORE metric and extended its framework to a multilingual setting. We conducted experiments using LLMs from the LLaMA, Qwen, and Aya families, generating biographies in 19 languages and comparing the results to Wikipedia pages. Our results reveal variations in hallucination rates, especially between high- and low-resource languages, raising important questions about LLM multilingual performance and the challenges in evaluating hallucinations in multilingual free-form text generation.

Keywords: LLM, Multilinguality, Hallucination, Factuality

1. Introduction

Since the public release of ChatGPT, large language models (LLMs) have gained popularity. They are increasingly being integrated into or even replacing traditional search engines, such as the LLaMA model for Meta mobile applications or Gemma for Google. This trend shows an increasing reliance on LLMs as sources of knowledge, due to their ability to generate human-like text. However, such use is concerning as LLMs tend to produce hallucinations.

A hallucination occurs when a LLM generates *false* content (Rawte et al., 2023) with respect to a specific *reference*. Based on the reference type, hallucinations can be classified as follows (Zhang et al., 2023c): input-conflicting, where the generated content contradicts the user’s input; context-conflicting, where it contradicts earlier outputs from the model; and fact-conflicting, where it contradicts established external knowledge. This work focuses exclusively on fact-conflicting hallucinations.

Understanding and tackling the issue of hallucinations in LLMs bring unique challenges. For example, detecting hallucinations is inherently difficult as they often appear highly credible. The wide range of tasks that LLMs are applied to also adds to the complexity, making it harder to comprehensively evaluate and mitigate hallucinations across different applications (Zhang et al., 2023c). Besides these well-known and investigated issues, hallucinations are also not produced in the same way depending on the prompt fed to the model. We introduce the concept of *multilingual hallucination gaps*, which refers to variations in the proportion of hallucinated outputs generated in response to prompts in different languages.

Measuring these gaps can reveal that prompts in certain languages are more likely to induce hallucinations than others, which can significantly impact the reliability and trustworthiness of LLMs, especially in low-resource languages.

Previous work (Hong et al., 2024; Lin et al., 2022) focused on measuring hallucinations through benchmarks that require human annotations, which can be costly and hard to scale for multilingual LLMs. These benchmarks are also not suited for a free-form text generation setup. As a result, an automated evaluation pipeline becomes highly desirable. Statistical measures like ROUGE fail to capture semantic variations (Sellam et al., 2020) while NLI-based approaches transfer poorly to these tasks (Falke et al., 2019). Among LLM-based methods, Chen et al., 2024 proposed an eigenvalue metric to measure self-consistent hallucination. However, measuring this metric is costly and might not be suitable for the evaluation of free-form generated text across multiple languages.

Consequently, we explore FACTSCORE (Min et al., 2023), a different LLM-based method to evaluate hallucinations. In particular, the FACTSCORE metric uses an LLM to fact-check outputs of other LLMs against a knowledge source. As the FACTSCORE was only developed and tested on English text, we extend the methodology to encompass different languages by comparing to knowledge sources in various languages and leveraging translation.

We evaluate LLMs from the LLaMA, Qwen and Aya families. We prompt them to generate biographies in 19 different languages and we then compute the FACTSCORE metric for each answer by comparing it to an external knowledge source, Wikipedia for this project. The computation is done through three different experimental setups. We finally analyzed the results with respect to the target language, to the experimental setup and to the LLM used for text generation. Our results show gaps in the FACTSCORE metric distribution across the prompt languages, particularly between high, medium, and low-resource languages. Our main contributions are:

1. Extending the FACTSCORE framework to a multilingual setting to quantify hallucination gaps across languages, with a focus on the disparities between high-resource and low-resource languages;

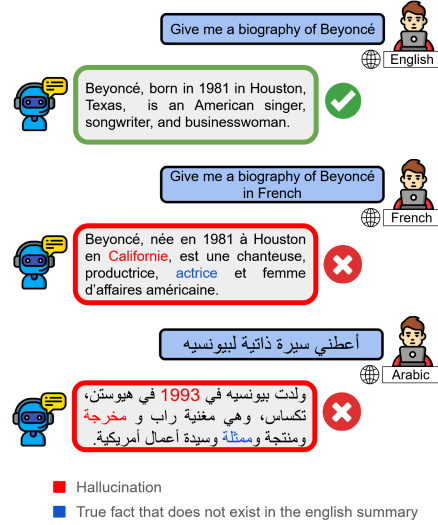


Figure 1: Example of Factual Hallucinations Gaps between languages

2. Evaluating a range of open-source and multilingual models to investigate improvements associated with different architectures and model sizes;
3. Assessing the robustness of the FACTSCORE framework across knowledge sources, prompt languages, and prompt templates.

2. Related work

Evaluating hallucinations Previous research has concentrated on evaluating, explaining, and mitigating hallucinations in language models (Ji et al., 2023; Zhang et al., 2023c). All these efforts have been focused on detecting hallucinations in English-generated text.

There are several human-annotated benchmarks available for this purpose, including those compiled in the unified benchmark on HuggingFace by Hong et al. (2024). Since these benchmarks rely on human annotation, they usually focus on short answers and are time-consuming to create, making them ill-suited for evaluating multilingual hallucinations in free-form text generation.

Automatic metrics for measuring hallucinations encompass statistical and model-based ones (Ji et al., 2023), many of which draw inspiration from summarization evaluation. Statistical metrics like ROUGE can only handle lexical information and fail to deal with syntactic or semantic variations (Sellam et al., 2020). NLI-based approaches are robust to lexical variability, but NLI models transfer poorly to abstractive summarization (Falke et al., 2019) and struggle to locate specific errors in generated content. Faithfulness Classification metrics (Liu et al., 2022) address this issue, but they rely heavily on English-annotated datasets.

Among LLM-based methods, Chen et al. (2024) proposed an eigenvalue-based metric for detecting self-consistent hallucinations. However, this approach is not well-suited for free-form text generation, where repeated prompts can produce different, yet correct, responses, leading to lower scores despite valid outputs. More relevant to long-form text generation are the methods proposed by Min et al. (2023) and Farquhar et al. (2024), both of which decompose answers into atomic facts. Min et al. (2023) employs a LLM to fact-check these facts against a knowledge source, while Farquhar et al. (2024) uses semantic entropy probabilities. In this study, we adopt the FACTSCORE (Min et al., 2023) approach, computationally less expensive.

Multilingual LLM Several studies have focused on evaluating language generation models within a multilingual framework. Some of these datasets include M3Exam (Zhang et al., 2023a) for performance on human exam and Flores-101 (Goyal et al., 2022) for translation abilities. For the performance of the LLMs we used on these datasets, refer to Annex B. We can note that these evaluation metrics are still not consistently disclosed in technical reports or widely-recognized benchmarks. Despite covering a broad range of applications, these datasets do not cover hallucinations. However, they do provide evidence that LLMs exhibit different performance across different languages.

An additional open research question concerns how multilingual abilities in these models are acquired (Zhang et al., 2023b; Wendler et al., 2024), as some models demonstrate proficiency in languages that are not officially supported. This observation motivated our decision to test models across a wide range of languages, even those not explicitly supported.

Hallucination metrics for multilingual generation Kang et al. (2024) examine automatic hallucination detection metrics across different languages, including ROUGE, Named Entity Overlap, and the NLI-based SUMMAC score. Their findings show that these metrics do not correlate. Previous studies have suggested that these metrics may not be reliable for assessing hallucinations (Ji et al., 2023), which motivates our investigation into LLM-based metrics, specifically the FACTSCORE metric, to evaluate hallucinations across languages on a range of open-source models.

The most related work to ours is Shafayat et al. (2024), as the authors also study how to extend the FACTSCORE metric to a multilingual context. However, their methodology revolves around prompting in the original language and then translating generated content before assessing factuality. We broaden our investigation by adding an experiment that prompts in English while requesting answers in another language, as well as another experiment that directly compares generations to the original language Wikipedia page. Further, we investigate the reliability of this choice of knowledge source. We also explore a wider range of languages, a different set of entities beyond politicians, as well as multilingual open-source models instead of ChatGPT. Additionally, our work critically examines the robustness of the metric itself and identifies areas for improving its reliability.

3. Measuring factuality

To evaluate factual hallucinations in multilingual free-form text generation, we use the FACTSCORE metric (Min et al., 2023). This metric is particularly suited for our goal because it offers an intuitive, automated evaluation pipeline that can be easily adapted to different languages. By breaking down responses into atomic facts, FACTSCORE not only provides a more precise measure of factuality but also provides two key pieces of information: the factuality rate and the number of facts in the response.

Let’s suppose we have a response \mathcal{R} generated by an LLM, hereinafter referred to as LM_{SUBJ} . The FACTSCORE metric for this response \mathcal{R} then consists of the following steps:

1. Decompose \mathcal{R} into a set of atomic facts $\mathcal{A}(\mathcal{R})$. An atomic fact is a short sentence conveying a single piece of information. This is achieved by prompting an LLM, hereinafter referred to as LM_{EVAL} , to "Please breakdown the following sentence into independent facts" after showing it some decomposition examples.
2. Compare each fact $a \in \mathcal{A}(\mathcal{R})$ with an external knowledge source \mathcal{C} . To do so, we retrieve the proper passage from the source \mathcal{C} . We then construct a prompt by concatenating the retrieved passage, the given atomic fact and "True or False?". We then feed this prompt to an LLM. We use the same LLM that was used to decompose atomic facts, the LM_{EVAL} . The answer gives us:

$$\text{Supported}(a, \mathcal{C}) = \mathbf{1}\{a \text{ supported by } \mathcal{C}\}$$

3. (Optional) Add a length penalty p depending on a hyperparameter γ if our response \mathcal{R} does not contain enough facts:

$$p = \exp\left(\frac{1 - \gamma}{|\mathcal{A}|}\right) \quad \text{if } |\mathcal{A}(\mathcal{R})| \leq \gamma$$

For instance, without applying this penalty, a response containing only one correct fact would receive a 100% FACTSCORE score, while a response with hundreds of facts, 99 of which are accurate, would get 99%. We set the default parameter to $\gamma = 10$, meaning responses with fewer than 10 facts are subject to a penalty.

4. Compute the FACTSCORE $F(\mathcal{R}, \mathcal{C})$:

$$F(\mathcal{R}, \mathcal{C}) = \frac{p}{|\mathcal{A}(\mathcal{R})|} \sum_{a \in \mathcal{A}(\mathcal{R})} \text{Supported}(a, \mathcal{C}) \quad (1)$$

4. Methodology

In this section, we explain and discuss the experimental settings, i.e. the choices of LM_{EVAL} , LM_{SUBJ} , content to be generated and knowledge source, then detail the experimental process.

4.1. Experimental settings

Experiments are built following the methodology of the FACTSCORE paper (Min et al., 2023). We prompt a LM_{SUBJ} to generate content in different languages and compute the FACTSCORE metric [1] with an LM_{EVAL} for each answer, by comparing it to an external knowledge source, specifically Wikipedia for this project.

Choice of the LM_{subj} Our objective is to evaluate open-source models with strong multilingual capabilities and developed in various countries. We also want to include for each model at least two different sizes to assess the impact of model size on the FACTSCORE. The final choice was set on LLaMA-3 (8B and 70B parameters), Aya-23 (8B and 35B parameters) and Qwen2 (7B and 72B parameters). Refer to Annex B for more details on the LLMs chosen. In the rest of the paper, we will refer to these models without specifying their versions.

Choice of the LM_{eval} To reduce bias in the evaluation process, we opt to use a different LM_{EVAL} than the LM_{SUBJ} models selected for the study. We use Mistral-7B-Instruct-v0.3 as the LM_{EVAL} and first compute the error rate and $F1_{\text{micro}}$ metrics (as detailed in Min et al., 2023) by computing FACTSCORE with Mistral on human annotated data. In Min et al., several methods are employed:

- No-context LM: Prompt ”<atomic-fact> True or False?”;
- Retrieve→LM: Retrieve passages from a knowledge source, concatenate these with the atomic fact and ”True or False?” and prompt the concatenated result;
- Nonparametric Probability (NP): Mask each token in the atomic fact, calculate its likelihood with a nonparametric masked LM, average probabilities, and make a prediction based on thresholding;
- Retrieve→LM + NP: Assign ”Supported” only if both Retrieve→LM and NP assign ”Supported”.

Based on their findings, we limit our experiments to the methods involving retrieval, comparing both Retrieve→Mistral and Retrieve→Mistral+NP. The results, presented in Table 3 in Annex C, show that Mistral performs competitively compared to the models used in

Min et al., validating our choice of Mistral as LM_{EVAL} . However, unlike with the model Inst-LLaMA, adding NP did not enhance performance. We will thus use the Retrieve→Mistral method for all subsequent experiments.

Prompts The FACTSCORE metric can be applied to any task, provided an appropriate knowledge source is available. We choose biographies because their factuality is easier to assess, as they generally include verifiable details such as birth dates and significant events, and they cover a wide range of nationalities. Besides, Wikipedia offers a multilingual knowledge source for this task, with biographies available in multiple languages.

We now have to choose the people whose biographies we will ask for and in which languages. For language selection, our goal is to cover a range of languages that includes high-resource, medium-resource, and low-resource languages. We define languages categories, i.e. high, medium and low, based on their data ratios from the Common Crawl corpus, drawing inspiration from Lai et al., 2023. For example, Spanish and Chinese are high-resource languages, Persian and Hindi fall under medium-resource, and Tamil and Swahili are examples of low-resource languages.

We also want to include as many language families as possible, while keeping the world’s most widely spoken languages. Since we will be using Wikipedia as the knowledge source, it is also important to ensure that all selected languages have sufficient Wikipedia coverage, which we measure by the number of Wikipedia pages available in each language. Table 4 in Annex D provides details on the 19 languages chosen, including the key statistics used to perform the selection. We do not take into account the languages supported by the LM_{SUBJ} in this selection, as some models demonstrate proficiency in languages that are not officially supported (Zhang et al., 2023b; Wendler et al., 2024). We then proceed to select a set of notable figures with Wikipedia pages available in all these languages, resulting in a list of 485 individuals. This selected set of entities is interestingly biased. Figures 5 and 6 in Annex D illustrate their top 15 countries of citizenship and languages spoken, respectively. The distribution is largely skewed towards the American citizenship and English language.

Wikipedia as a knowledge source We retrieve Wikipedia summaries in every language for the 485 notable figures to serve as the knowledge source. Recognizing that Wikipedia content quality can vary across languages, we start by comparing these summaries to each other. For a given language, we translate the Wikipedia pages with GPT-4o-mini and we compute the FACTSCORE comparing to the English Wikipedia. We choose the English Wikipedia because the FACTSCORE indicates whether an atomic fact is supported by the knowledge source, rather than giving information on its presence in the knowledge source. Our assumption is that English Wikipedia is more comprehensive than other language versions. We also compute FACTSCORE for the English Wikipedia pages, i.e. comparing them against themselves. Results are presented in Figure 8 in Annex E.

We can directly see that the evaluator is not perfect. Indeed, comparing the English Wikipedia to itself does not always yield a 100% FACTSCORE, even if it typically falls within a high range of 90-100%. For the other languages, cross-checking with English yields much lower FACTSCORE. This suggests that content in different languages can either contradict or diverge from what is found in English Wikipedia. These observations are important for the rest of our analysis. We will compare a generated biography with both its original language Wikipedia version and the English one. This may give us insight on how the

LLM captures and represents knowledge in a multilingual setting, depending on whether the FACTSCORE is different when using these two different knowledge sources. It is also important to consider both for a more precise analysis.

4.2. Experiments

In this section, we outline the experimental process, detailing the steps from data generation to FACTSCORE evaluation, with intermediate sanity checks.

Data generation To ensure robustness in measuring the hallucination gaps, we use three established prompt templates for generating biographies from the literature: "Tell me a biography of {}", "Give me a biography of {}" and "Please give me a biography of {}".

We use two prompting methods:

- **lang-prompt**: Translate the template in the target language **lang** and use the translated prompt;
- **en-prompt**: Use the English template but add "in {lang}" at the end (e.g., "Tell me a biography of {} in French").

We use these methods for our six models LM_{SUBJ} , resulting in a total of 12 text generation setups. Each setup produces $485 \times 19 \times 3$ generated responses.

Sanity checks Once the biographies are generated, we carry out some sanity checks to ensure we have good enough generated answers to compute FACTSCORE with. For each answer, we verify that it is in the correct target language with the module `py3langid`. We set a threshold of 20 distinct words to remove outputs with same words repeated infinitely. We do not compute FACTSCORE for generated answers that do not pass sanity checks.

FactScore evaluation We perform three experiments depending on the knowledge source and prompting method used. We will refer to these experiments as (prompt language, Wikipedia language):

1. (**lang**, **lang**): Compare the response produced with a **lang**-prompt to the **lang** Wikipedia page;
2. (**lang**, **en**): Translate the response produced with a **lang**-prompt to English and compare to the English Wikipedia page;
3. (**en**, **en**): Translate the response produced with an **en**-prompt to English and compare to the English Wikipedia page.

Recall that we generated three answers for each notable figure and language. We thus compute the FACTSCORE score for each answer and then average these three scores to obtain a single final score per entity. We also report the standard deviation. This process is repeated across all three experiments: (**lang**, **lang**), (**lang**, **en**) and (**en**, **en**).

Translation All translation steps, including prompt translations, generations, and Wikipedia pages, were performed with GPT-4.

5. Results

5.1. Sanity checks

The percentages of generated answers that passed sanity checks for each LM_{SUBJ} are shown in Table 5 in Annex F. These initial results already demonstrate whether the LM_{SUBJ} is able to generate multilingual text, including in languages stated as not supported. Figure 2 shows the percentage of generated responses produced in the correct target language, for each LM_{SUBJ} and prompt setting. As expected, the models generally perform better in generating text in high-resource languages compared to low- and very low-resource ones. For very low-resource languages like Javanese (jv) and Malay (ms), the models rarely generate text in the correct target language.

The best performance overall is achieved with Qwen7 both in English (**en**) and original language (**lang**) prompting. We can note that for all models prompting in English tend to decrease the percentages of generated responses in the correct target language. Interestingly, for the LLaMA and Qwen families of models, increasing the number of parameters does not always lead to better performance, especially when prompted in English.

We also observe poorer performance in the Japanese, Chinese, and Korean languages for the LLaMA models compared to others. When examining generated answers that failed the sanity checks, we notice that the LLaMA models often produced Romaji (writing in Roman characters) instead of Kanji (writing using Chinese characters).

5.2. FACTScore

The LM_{subj} studied show different trends in hallucination gaps. Figure 3 illustrates the mean FACTScore per model and per language for the (**en**, **en**) experiment. For results from the two other experiments, see Annex G. The model Qwen72 performs the best overall. This aligns with the models’ multilingual performance on other benchmarks, where Qwen72 also ranks highest (see Table 2).

Studying different subject models shows the impact of training sets, architectures, and model sizes on the FACTScore results. Aya and Qwen models, with extensive multilingual training data, have good performance in supported languages (see Annex B), outperforming LLaMA models, which are not optimized for multilingual tasks. Notably, Qwen models achieve the highest performance in Chinese, benefiting from tokenizers specifically designed for Chinese characters. Furthermore, within each model family, a higher number of parameters correlates with higher FACTScore. These findings reveal that multilingual hallucination gaps are influenced not only by high- versus low-resource language distinctions but also by factors like tokenizer design and language distribution in training data.

Hallucination rates differ across languages. Table 1 presents the average results across all models and entities, grouped by language category. We observe multilingual hallucination gaps, with both factuality and the number of facts decreasing as the language resource level decline. Due to high standard deviations within these categories, we examine the FACTScore distribution at a finer level, for each language.

Figure 4 shows these distributions across all models and entities, broken down by language and experiment. See Annex I for the same results per LM_{SUBJ} . Aside from Malay (ms) and Javanese (jv) in all experiments, and Japanese (ja) and Chinese (zh) languages in the

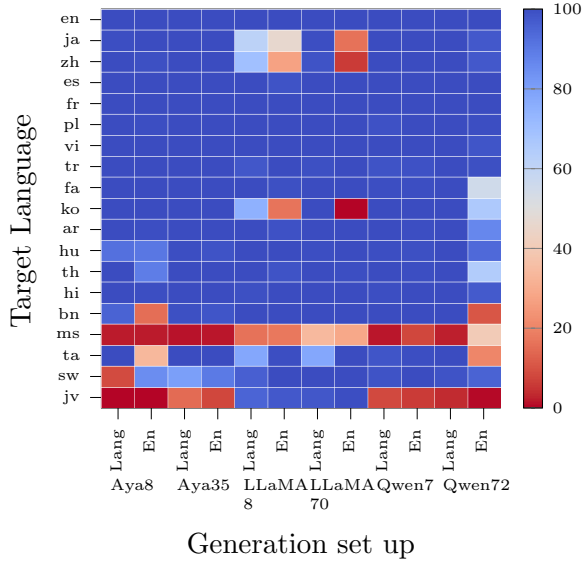


Figure 2: Percentage of correct language produced per target language and generation set up

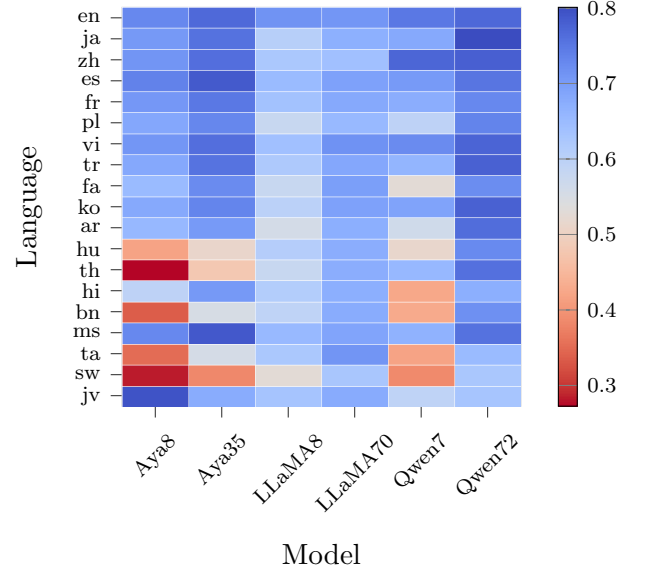


Figure 3: FACTScore per language and per model for the (en, en) experiment

Language Category	FACTScore (%)			# of Facts		
	(en, en)	(lang, en)	(lang, lang)	(en, en)	(lang, en)	(lang, lang)
Very-High	73.7 (\pm 10.1)	71.8 (\pm 9.9)	70.3 (\pm 9.7)	79	82	103
High	70.2 (\pm 12.6)	69.3 (\pm 13.4)	58.5 (\pm 16.4)	68	73	65
Medium	64.7 (\pm 16.0)	61.3 (\pm 19.5)	47.8 (\pm 19.4)	54	59	49
Low	56.9 (\pm 18.9)	47.6 (\pm 23.4)	44.4 (\pm 20.4)	38	34	53

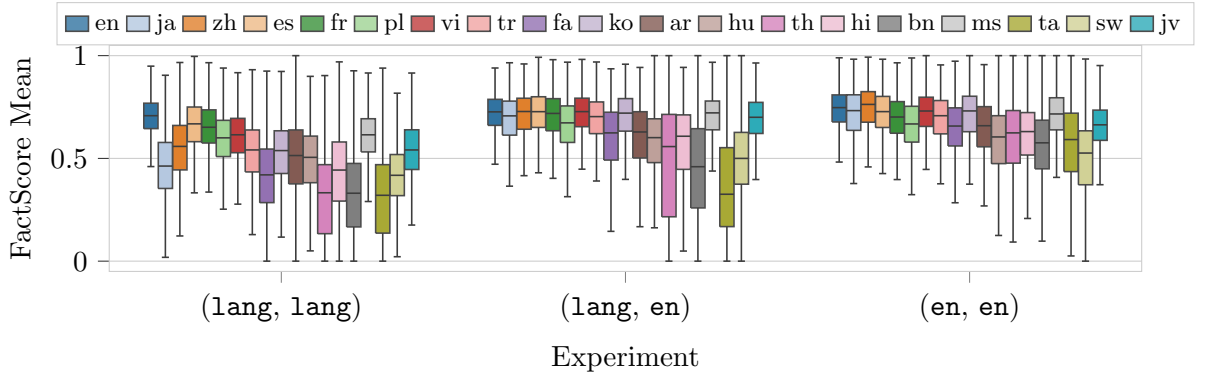
 Table 1: Mean FACTScore (\pm STD) and Mean number of facts by Language Category and Experiment for all models


Figure 4: FactScore Mean distribution by Language and Experiment for all models

(**lang**, **lang**) experiment, we can observe the same trend as FACTSCORE distributions are more spread out and shifted toward lower values as the language resource level decline. It should be noted that after filtering out unsane answers that were in the incorrect language, we have less data points for the Malay and Javanese languages (see Figure 2). On the contrary, only for the (**lang**, **lang**) experiment, Japanese and Chinese show distributions that are more spread out and shifted toward lower values compared to other high-resource languages. This raises the question of why different experimental setups affect the results.

Different pipelines show different results. Table 1 and Figure 4 highlight differences in results across the (**lang**, **lang**), (**lang**, **en**), and (**en**, **en**) experiments, showing that the choice of knowledge source and prompt language can influence the FACTSCORE outcomes.

Regarding the prompt language, we gain insights into the performance of the LM_{SUBJ}. As shown in Figure 2, models respond differently to prompts given in English versus the target language **lang**. When comparing (**lang**, **en**) and (**en**, **en**) — where the knowledge source remains the same but the prompt languages differ — we observe the highest results with the (**en**, **en**) setting. This raises concerns, as we would want models to respond accurately to prompts in original languages.

The impact of the knowledge source presents more significant challenges, as it directly affects the evaluator. In comparing (**lang**, **lang**) and (**lang**, **en**), where the prompt language remains the same but the knowledge source differs, we see a decrease in FACTSCORE for the (**lang**, **lang**) experiment, and more dispersed distributions. This trend may be attributed to the quality of Wikipedia pages in their original languages, as highlighted in Figure 8.

The performance of (**lang**, **lang**) compared to (**lang**, **en**) and (**en**, **en**) is also influenced by the respective multilingual capabilities of the LM_{EVAL} and the translator (GPT-4). The (**lang**, **lang**) experiment is the only one that involves prompting the LM_{EVAL} in languages other than English. We noticed that for most languages, while breaking down a generation into atomic facts, the LM_{EVAL} also translated these facts into English even when not instructed to do so. Addressing this behavior might improve results in the (**lang**, **lang**) setting. The (**lang**, **en**) approach seems to be the most suitable option, as it allows us to maintain prompts in their original languages.

FactScore’s robustness depends on the language. Table 6 in Annex H presents FACTSCORE standard deviations across three prompt templates per entity. Averaging these by language category and experiment reveals that lower-resource languages exhibit higher FACTSCORE variability, indicating less consistency across prompt templates.

6. Conclusion

Our research shows multilingual hallucination gaps in LLMs: models tend to hallucinate more in low-resource languages and show greater factual accuracy in high-resource ones. Larger models generally perform better but still show hallucinations in low-resource languages. Even models with strong multilingual capabilities hallucinate in such languages, and struggle to generalize effectively to unsupported languages, suggesting that simply increasing model size or expanding training data is not a complete solution. This raises important

concerns about the equitable performance of LLMs across different linguistic groups and the broader implications for fairness in AI technologies.

Code availability

The code and data are accessible at the GitHub repository: https://github.com/cleachataig/Multilingual_Hallucination_Gaps.

Acknowledgments

Funding support for project activities has been partially provided by Canada CIFAR AI Chair, Facebook award, MEI award and FRQNT award. We also express our gratitude to Compute Canada and to Mila (mila.quebec) for their support in providing facilities for our evaluations.

References

- Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, Kelly Marchisio, Max Bartolo, Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. Aya 23: Open weight releases to further multilingual progress, 2024. URL <https://arxiv.org/abs/2405.15032>.
- Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. Inside: Llms’ internal states retain the power of hallucination detection. *arXiv preprint arXiv:2402.03744*, 2024.
- Abhimanyu Dubey et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Tobias Falke, Leonardo F. R. Ribeiro, Prasetya Ajie Utama, Ido Dagan, and Iryna Gurevych. Ranking generated summaries by correctness: An interesting but challenging application for natural language inference. In Anna Korhonen, David Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2214–2220, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1213. URL <https://aclanthology.org/P19-1213>.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630, 2024.
- Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. The flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538, 2022.

- Giwon Hong, Aryo Pradipta Gema, Rohit Saxena, Xiaotang Du, Ping Nie, Yu Zhao, Laura Perez-Beltrachini, Max Ryabinin, Xuanli He, and Pasquale Minervini. The hallucinations leaderboard—an open effort to measure hallucinations in large language models. *arXiv preprint arXiv:2404.05904*, 2024.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.
- Haoqiang Kang, Terra Blevins, and Luke Zettlemoyer. Comparing hallucination detection metrics for multilingual generation, 2024.
- Viet Lai, Nghia Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Nguyen. ChatGPT beyond English: Towards a comprehensive evaluation of large language models in multilingual learning. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13171–13189, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.878. URL <https://aclanthology.org/2023.findings-emnlp.878>.
- Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human falsehoods. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.229. URL <https://aclanthology.org/2022.acl-long.229>.
- Tianyu Liu, Yizhe Zhang, Chris Brockett, Yi Mao, Zhifang Sui, Weizhu Chen, and Bill Dolan. A token-level reference-free hallucination detection benchmark for free-form text generation. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6723–6737, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.464. URL <https://aclanthology.org/2022.acl-long.464>.
- Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.741. URL <https://aclanthology.org/2023.emnlp-main.741>.
- Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation models, 2023.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. BLEURT: Learning robust metrics for text generation. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault,

- editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.704. URL <https://aclanthology.org/2020.acl-main.704>.
- Sheikh Shafayat, Eunsu Kim, Juhyun Oh, and Alice Oh. Multi-fact: Assessing multilingual llms’ multi-regional knowledge using factscore, 2024.
- Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. Do llamas work in english? on the latent language of multilingual transformers. *arXiv preprint arXiv:2402.10588*, 2024.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- Wenxuan Zhang, Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models. *Advances in Neural Information Processing Systems*, 36:5484–5505, 2023a.
- Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. Don’t trust ChatGPT when your question is not in English: A study of multilingual abilities and types of LLMs. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7915–7927, Singapore, December 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.491. URL <https://aclanthology.org/2023.emnlp-main.491>.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren’s song in the ai ocean: A survey on hallucination in large language models, 2023c.