## Y-NQ:

# English-Yorùbá Evaluation dataset for Open-Book Reading Comprehension and Text Generation

Marta R. Costa-jussà, Joy Chen, Ifeoluwanimi Adebara, Joe Chuang, Christophe Ropers, Eduardo Sánchez

FAIR at Meta

The purpose of this work is to share an English-Yorùbá evaluation dataset for open-book reading comprehension and text generation to assess the performance of models both in a high- and a low-resource language. The dataset contains 358 questions and answers on 338 English documents and 208 Yorùbá documents. The average document length is  $\approx 10 \mathrm{k}$  words for English and 430 words for Yorùbá. Experiments show a consistent disparity in performance between the two languages, with Yorùbá falling behind English for automatic metrics even if documents are much shorter for this language. For a small set of documents with comparable length, performance of Yorùbá drops by x2.5 times. When analyzing performance by length, we observe that Yorùbá decreases performance dramatically for documents that reach 1500 words while English performance is barely affected at that length. Our dataset opens the door to showcasing if English LLM reading comprehension capabilities extend to Yorùbá, which for the evaluated LLMs is not the case.

Date: December 12, 2024

Correspondence: Marta R. Costa-jussà at costajussa@meta.com

Meta

## 1 Introduction

This study explores the intersection of reading comprehension and text generation, examining how models perform on tasks requiring both in-context understanding (i.e., open-book model, where the model has access to the context document during inference to answer a particular question) and generative text production (i.e. the answer is free-text which has to be compared to a gold standard reference). We aim to investigate the performance of this task in two languages: a high-resource language (English) and a low-resource language (Yorùbá). For this, we introduce Y-NQ (Yorùbá Natural Questions) a comprehensive open-book question-answer dataset (Section 2). Y-NQ is sourced from NQ (Kwiatkowski et al., 2019) and provides a complete article context for informed answers and text generation tasks, and parallel documents on the same topic for both high- and low-resource languages. The data set also includes the comparability of the responses in languages. As a result, we are increasing Natural Language Processing (NLP) resources in Yorùbá (Ahia et al., 2024). Our data set is benchmarked against state-of-the-art Large Language Models (LLMs). The results and analysis (Section 3) shows that responses in Yorùbá are more inaccurate than those in English.

As a by-product of human annotations, we identify inaccuracies in the English-language version of some Wikipedia articles (26 incorrect answers out of 1,566 humanly analyzed questions in the English-language subset of articles), which confirms the existence of accuracy discrepancies across languages for the same Wikipedia topics, thus supporting, for example, the need to better interlink Wikipedia articles across languages (Klang and Nugues, 2016).

## 2 Dataset description

#### 2.1 Requirements and Background

The performance of Reading Comprehension (RC) in LLMs has been explored in different settings. At the high level, RC tasks can fall under two main categories: open-book tasks, such as in SQuAD (Rajpurkar et al., 2016), and close-book tasks, such as in TriviaQA (Joshi et al., 2017). Response formats vary across RC tasks as well and include: true/false classification (e.g., BoolQ; Clark et al., 2019), multiple-choice questions (e.g., Belebele), span selection (e.g., SQuAD), and text generation (e.g., NQ or TriviaQA).

Since we are interested in exploring the intersection of reading comprehension and text generation covering both a high- and a low-resource language, we can explicitly set our requirements to include for each of the two types of language: (a) long articles (>100s words), (b) question-answer pairs with lengthy answers (>10s words), and (c) equivalence annotations for cross-lingual answers. Since there are no existing data sets to this effect, we extend existing research by tailoring an established data set to our specific requirements. We justify our choice of data sets and low-resource language selection as explained in the following.

**Dataset.** Among the open-book and text generation tasks, one of the largest datasets with multilingual information available is NQ.

Objective	Read an article and find a paragraph containing enough information to answer a specific question.
Project Context	Evaluate accuracy of large language models in finding long contexts and short answers; extend Natural Questions dataset to multilingual, non-English centric.
Task Components	<ul> <li>QUESTION: Simple question requesting information or explanation.</li> <li>ARTICLE: Numbered paragraphs containing relevant information.</li> </ul>
Task Steps	<ol> <li>Read QUESTION carefully.</li> <li>Read ARTICLE paragraphs until sufficient information is found.</li> <li>Record findings by answering task questions.</li> </ol>
Additional task steps	Discard questions that contain the answer in English in the Yorùbá document When possible, add Yorùbá questions, translate them into English, and find answers both in the Yorùbá and English documents.

**Table 1** Linguistic guidelines and annotation

Low-resource language. There is a large number of low-resource languages that could be explored here. We prioritize a low-resource language that has overall limited digital resources (in compliance with the definition of low resource), but has a high representation in Wikipedia (on the order of several thousands of entries) and a significant number of speakers (in the order of tens of millions), and makes use of the same script (Latin) as the high-resource language in which results are compared. One of the languages that complies with all these criteria is Yorùbá, in which we can also find works on comprehension of the language in the domain of language exams (Aremu et al., 2024), based on short passages and multiple choice answers. Another work is the AfriQA dataset (Ogundepo et al., 2023) for answering open-retrieval questions, with a primary focus on retrieving correct answers that are answerable on Wikipedia. However, this cannot be used as an open book. Finally, Bebebele (Bandarkar et al., 2024) also includes Yorùbá, although it uses short passages and multiple choice answers.

#### 2.2 Dataset creation

**NQ pre-selection.** We looked at 315,203 examples and 231,695 unique English Wikipedia pages from the NQ training and validation datasets. We filter questions for only those where every long answer is contained in an html tag where is the first identified html tag in the long answer span. This filters out about 25 percent of the questions.

We extracted 2,855 Yorùbá Wikipedia pages that are actively associated with the above English pages. We removed documents with fewer than 500 characters, including formatting, and performed multiple cleaning procedures, such as removing html formatting, removing citation notations, and filtering out irrelevant sections in Wikipedia articles (e.g., references, tables). 664 Yorùbá documents and 1,566 questions were sent for human annotation.

**Pre-annotation effort.** In order to reduce the annotation workload, we automatically pre-selected Yorùbá sentences that could be good response candidates by computing a similarity score. If the answer to the question was in agreement with a high similarity score, the annotator would save time by looking through the document and only checking if the match was correct. We conducted a SONAR embedding similarity (Duquenne et al., 2023) analysis between Yorùbá documents and long English answers. We used the Stopes sensitizers on all text extracted from elements for both the scraped Yorùbá Wikipedia articles downloaded from the previous step and the original NQ Wikipedia pages. We then created SONAR embeddings of each extracted sentence and identified those sentences in the Yorùbá pages which were most similar to sentences in the long English answers based on their cosine similarity scores. For a small set of samples, we asked the annotators to examine the entries in a small validation data set to identify a reasonable threshold indicating high similarity between Yorùbá/English sentences, which could then be applied to the rest of the data set. The analysis shows a low similarity matching rate, which is likely due to the low quality and short length of many Yorùbá articles and/or SONAR embeddings not being suitable for such a task. Given this low reliability, we abandoned this automatic pre-annotation, which would not reduce annotation efforts.

**Annotation guidelines and requirements.** We designed the annotation guidelines as follows. We provided context on the objective of the task together with the project context and description of the task. The guidelines are summarized in Table 1.

Finally, beyond the guidelines, we provided additional examples and requested that annotators should be native speakers of the language of the source documents and should have at least CEFR C2 level proficiency in English.

	Eng	Yor
#Q&A	358	358
#DOCS	338	208
AVG. DOC LEN	10363	430
MEDIAN DOC LEN	9272	172
AVG. QUESTION LEN	8.86	9.39
AVG. LONG ANSWER LEN	113.80	32.89

**Table 2** Dataset Statistics. Length is in words.

**Annotator findings.** We noticed that many articles have a significant amount of English content. Several documents also contained errors, such as incorrect spelling, ungrammatical sentences, and sentences that lacked clarity or meaning. We disregarded such articles and corrected articles that were contaminated with a

<sup>&</sup>lt;sup>1</sup>https://github.com/facebookresearch/stopes

small amount of English content. We also removed the entries where no answers could be found in the Yorùbá articles.

Following the guidelines, the annotators encountered the following: (a) questions with multiple correct answers, for which they annotated each correct answer for the question; (b) questions with correct answers in Yorùbá, but incorrect in English, where they annotated the Yorùbá appropriately, but flagged the English portion incorrect (there were 26 questions in the category); (c) unclear questions (5 questions) to which no annotations were assigned; (d) answers existing in multiple paragraphs in the document for which they annotated the row with all paragraphs where

There were 456 Yorùbá documents that did not answer the question; therefore, we discarded those. Only eight incorrect English answers from the previous 26 remain in the final dataset, and we did not correct them since the English documents remained the same as in the original NQ.

**Statistics.** Table 2 details the statistics of the data set<sup>2</sup>. Our carefully curated selection contains 208 unique Yorùbá Wikipedia documents with an average word count of 430, and 356 unique questions. Only the questions are strictly comparable. English and Yorùbá documents are not comparable in number or length, but they are so in topic and domain. The answers are not comparable in length. Notice that English documents outnumber Yorùbá documents mainly due to multiple versions of the same English topic counted as different documents, while in Yorùbá we selected one version of the document and multiple topics in English that correspond to the same Yorùbá topic.

The fact that English documents are longer than those in Yorùbá makes the task easier for Yorùbá, since documents are significantly shorter within the same topic or domain. We identified a subset of six documents that are strictly comparable in length and topic for English and Yorùbá, which allows us to make a fair comparison. Table 3 shows the list of fields in Y-NQ and a sample entry.

FIELD	DESCRIPTION	Example
1. Question ID	Unique identifier	3506772758530306034
2. English Document	English text document	
3. English Question	Question in English	what is the name of the first nigerian
		president
4. English Long Answer	Detailed answer in English	.ky is the Internet country code top-level
		domain (ccTLD) for the Cayman []
5. English Short Answer	Brief answer in English	Nnamdi Azikiwe
6. Yorùbá Document	Yorùbá text document	
7. Yorùbá Rewrite Flag	Was Yorùbá document rewritten?	1
	(0: no, 1: yes)	
8. Yorùbá Question	Question in Yorùbá	kí ni ky dúró fún ní erékùṣù cayman
9. Yorùbá Short Answer	Brief answer in Yorùbá	Nnamdi Azikiwe ni Aare
10. Yorùbá Long Answer	Detailed answer in Yorùbá	Nnamdi Azikiwe ti o je Gomina Agba
		nigbana di Aare, ipo to je fun ayeye, []
11. Yorùbá Paragraph Info	Contextual information	P2
12. Answer Alignment	Semantic equivalence	1
	(0: not literal, 1: literal)	

**Table 3** Dataset Fields, Descriptions and Sample entry.

<sup>&</sup>lt;sup>2</sup>There are two questions that come from the validation NQ dataset, which have two different answers

	Lan	R-1	R-2	R-L
GPT40	Eng	0.39	0.23	0.30
	Yor	0.34	0.19	0.27
O1MINI	Eng	0.45	0.22	0.30
	Yor	0.30	0.14	0.22
LLAMA	Eng	0.31	0.18	0.23
	Yor	0.20	0.15	0.18

**Table 4** Results for 3 LLM in terms of Rouge computed for the entire set of questions. Human Score is computed on 358 questions.

## 3 Experiments

**Baselines** We evaluate our dataset with GPT-4o<sup>3</sup> (et al., 2024b), o1-mini<sup>4</sup>, and LlaMA-3.1-8b (et al., 2024a), therevy covering both open and closed models, as well as models of different sizes. For each Y-NQ entry, we prompt the models with the following formatted instructions.

,, ,, ,,

Given the following passage and a question, answer the question in a single paragraph with information found in the passage.

#### PASSAGE

{document}
####
OUESTION

{question} ####

**ANSWER** 

,, ,, ,,

**Evaluation.** We evaluate the results by comparing the generated text and the reference long answer using several Rouge (Lin, 2004) versions (Rouge-1, Rouge-2, Rouge-L).

**Automatic metrics.** Table 4 reports the results showing that Yorùbá consistently performs worse than English (e.g., losing 0.4 in Rouge-1). However, the Yorùbá task is much easier because the documents are much shorter, which means that answering the question becomes an easier task. Even if we prompt the model to only answer based on the in-context document, we can not discard the idea that English may get better results due to using the internal knowledge from the model.

**Length analysis.** Model performance changes with the length of the document, as shown in Figure 1. The dataset was split into equal size of documents in each length bucket. We can see a drop in performance

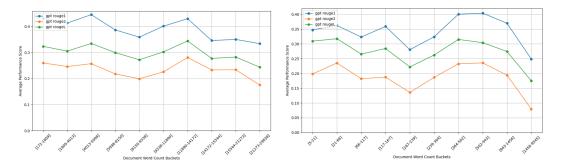
<sup>&</sup>lt;sup>3</sup>gpt-4o version 2024-08-06

<sup>&</sup>lt;sup>4</sup>o1-mini version 2024-09-12

when the Yorùbá documents reach 1,500 words, which shows the challenges that current models face in long-context understanding of low-resource languages. For a small portion of long-enough documents of comparable length between English and Yorùbá (only 4 documents that are over 900 words long), English performance demonstrates a significant edge (1.58X-2.56X), see Table 5.

	AVG W.	R-1	R-2	R-L
Eng	3299	0.45	0.23	0.30
Yor	3070	0.32	0.09	0.19

**Table 5** Results for six comparable English and Yorùbá documents



**Figure 1** Impact of Document Length Buckets on Performance Scores for English (top) and Yorùbá (bottom) for GPT-4 outputs

#### 4 Conclusions

Y-NQ is a newly released dataset that enables to compare generative open-book reading comprehension between English and Yorùbá. The main contributions of our data set are to allow for the comparison of LLM results in a reading comprehension task across a high- and a low-resource language, showing what are the generalization capabilities of LLMs in this particular case. Moreover, our annotations confirmed variations in the accuracy of Wikipedia articles in all languages. In particular, we identify inaccurate English responses for Yorùbá language-specific content. Y-NQ allows us to evaluate how reading comprehension capabilities extend to Yorùbá. Y-NQ is not exactly comparable in its totality between languages. Given that Yorùbá has shorter documents than English, the reading comprehension task is easier for Yorùbá. Therefore, results on this language should be much better than in English to expect parity between languages. Our experiments show that the reading comprehension capabilities of current English LLMs do not extend to Yorùbá. Y-NQ is freely available on HuggingFace.

## **Limitations and Ethical considerations**

Y-NQ is limited in size, language, and domain coverage. The fact of using Wikipedia and extending an existing open-source dataset (NQ) may play in favor of having higher results in both languages due to contamination. Furthermore, the data set is not fully comparable between English and Yorùbá, since documents and answers vary in length.

Our experimentation is limited to models and automatic evaluation metrics, which could be compensated for through human evaluation. Annotators were paid a fair rate.

## **Acknowledgements**

This paper is part of the LCM project<sup>5</sup> and the authors would like to thank the entire LCM team for the fruitful discussions.

### References

Orevaoghene Ahia, Anuoluwapo Aremu, Diana Abagyan, Hila Gonen, David Ifeoluwa Adelani, Daud Abolade, Noah A. Smith, and Yulia Tsvetkov. Voices unheard: NLP resources and models for Yorùbá regional dialects. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4392–4409, Miami, Florida, USA, November 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.emnlp-main.251.

Anuoluwapo Aremu, Jesujoba O. Alabi, Daud Abolade, Nkechinyere F. Aguobi, Shamsuddeen Hassan Muhammad, and David Ifeoluwa Adelani. Naijarc: A multi-choice reading comprehension dataset for nigerian languages, 2024. URL https://arxiv.org/abs/2308.09768.

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. The belebele benchmark: a parallel reading comprehension dataset in 122 language variants. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 749–775, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.44. URL https://aclanthology.org/2024.acl-long.44.

Christopher Clark, Matthew Gardner, Tom Fevry, and Robert Weischedel. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv* preprint arXiv:1905.10044, 2019.

Paul-Ambroise Duquenne, Holger Schwenk, and Benoît Sagot. Sonar: Sentence-level multimodal and language-agnostic representations, 2023. URL https://arxiv.org/abs/2308.11466.

Abhimanyu Dubey et al. The llama 3 herd of models, 2024a. URL https://arxiv.org/abs/2407.21783.

OpenAI et al. Gpt-4 technical report, 2024b. URL https://arxiv.org/abs/2303.08774.

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. Triviaqa: A large scale question-answer dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.

Marcus Klang and Pierre Nugues. Pairing Wikipedia articles across languages. In Key-Sun Choi, Christina Unger, Piek Vossen, Jin-Dong Kim, Noriko Kando, and Axel-Cyrille Ngonga Ngomo, editors, *Proceedings of the Open Knowledge Base and Question Answering Workshop (OKBQA 2016)*, pages 72–76, Osaka, Japan, December 2016. The COLING 2016 Organizing Committee. URL https://aclanthology.org/W16-4410.

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. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*, 2019.

Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://aclanthology.org/W04-1013.

Odunayo Ogundepo, Tajuddeen Gwadabe, Clara Rivera, Jonathan Clark, Sebastian Ruder, David Adelani, Bonaventure Dossou, Abdou Diop, Claytone Sikasote, Gilles Hacheme, Happy Buzaaba, Ignatius Ezeani, Rooweither Mabuya, Salomey Osei, Chris Emezue, Albert Kahira, Shamsuddeen Muhammad, Akintunde Oladipo, Abraham Owodunni, Atnafu Tonja, Iyanuoluwa Shode, Akari Asai, Anuoluwapo Aremu, Ayodele Awokoya, Bernard Opoku, Chiamaka Chukwuneke, Christine Mwase, Clemencia Siro, Stephen Arthur, Tunde Ajayi, Verrah Otiende, Andre Rubungo,

<sup>&</sup>lt;sup>5</sup>https://github.com/facebookresearch/large\_concept\_models

Boyd Sinkala, Daniel Ajisafe, Emeka Onwuegbuzia, Falalu Lawan, Ibrahim Ahmad, Jesujoba Alabi, Chinedu Mbonu, Mofetoluwa Adeyemi, Mofya Phiri, Orevaoghene Ahia, Ruqayya Iro, and Sonia Adhiambo. Cross-lingual open-retrieval question answering for African languages. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14957–14972, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.997. URL https://aclanthology.org/2023.findings-emnlp.997.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.