

EIFFEL: a novel benchmark to measure bias of English heavy training on French idiomatic expressions

Anonymous ACL submission

Abstract

Mainstream multilingual LLMs are generally trained on a much higher proportion of English than multilingual data, raising questions about their ability to capture linguistic features particular to non-English languages or to capture information important to non-anglophone cultures. We add to a growing effort to increase multilingual sensitivity in LLMs by developing a benchmark, EIFFEL, testing mastery of French idiomatic expressions in context. We fully explain the methodology, which exploits input from native French speakers, to make it reproducible for other languages. We compare mainstream multilingual LLMs with French-focused LLMs both on standard LLM benchmarks and EIFFEL; EIFFEL brings out the benefits of higher proportions of French data and shows limitations of standard benchmarks for measuring multilingual competence. We also train from scratch a series of 1B SLMs with different proportions of French and English pre-training data that confirm EIFFEL’s lessons.

1 Introduction

While large language models (LLMs) are increasingly popular worldwide, many of the leading models are trained on disproportionate amounts of English data. For example, only 8% of Llama 3.1’s training data come from non-English natural languages (Grattafiori et al., 2024). This raises the question of how anglocentricity shapes an LLM’s ability to produce high-quality sequences in other languages and to represent knowledge and cultural norms central to non-anglophone cultures.

Answering this question is complicated by the potential for language transfer. Suppose that we have a bilingual model B trained on English and some non-English language L . If there is transfer from one language to another, then B ’s probability distribution over English tokens can inform its distribution over L tokens and vice versa. If we then apply B to a downstream task that is covered

by knowledge transfer from English, B might do well having seen only a small proportion of data in L . Language transfer together with task relevant training is arguably what makes mainstream anglocentric LLMs surprisingly powerful in non-English languages. A more focused question is then: to what extent will culturally and linguistically sensitive aspects of L be overlooked if we rely on language transfer?

A complete answer to this question is beyond the scope of this paper, but we offer an important and first of its kind tool to explore this question: EIFFEL, Evaluation of Idiomatic French Fixed Expressions for Large Language Models, is a benchmark that showcases French idiomatic expressions.

Idiomatic expressions make a good test subject because they are a feature of everyday language that is highly language specific. While some idiomatic expressions can be translated word-for-word between French and English, such as “Not my cup of tea/Pas ma tasse de thé,” others are less direct or even completely different. The expression “call a spade a spade” for example, has an obvious counterpart in French, literally “call a cat a cat”; but it is not a direct translation. The majority of expressions are even less easily translatable. “Avoir du chien,” literally, “to have some dog,” means to be charming. For a model to handle these latter expressions, we hypothesize that it needs more than a solid hold on English and a good capacity for translation; it needs to have seen either explicit translations of the idiomatic expressions or a fair amount of (non-translated) French data.

We test the impact of our benchmark by first comparing a series of models on French and English versions of standard benchmarks (ARC Challenge (Clark et al., 2018), Hellaswag (Zellers et al., 2019), MMLU (Hendrycks et al., 2020)) and then compare the results to those on EIFFEL and the French subset of INCLUDE (Romanou

et al., 2024), a dataset also originally in French that contains culturally sensitive and agnostic subsets that allow for interesting comparisons with EIFFEL. We consider pretrained models from two anglocentric model families,¹ Llama and Gemma (Team et al., 2024, 2025), as well as “gallo-centric” models² trained on 1:2 and 1:1 ratios of French and English data: the Gaperon models (Godey et al., 2025), Lucie (Gouvert et al., 2025), and CroissantLLM (Faysse et al., 2024). We also look at EuroLLM models (Martins et al., 2025a,b), as they offer a middle case between anglocentric and gallo-centric models and focus on translation capacities. Our study is restricted to pretrained models, as we are interested in basic linguistic mastery.

While standard LLM benchmarks in French do not reveal an advantage for gallo-centric models over anglocentric ones, both EIFFEL and the culturally-sensitive subset of INCLUDE do. This suggests that EIFFEL captures features of French that benefit less from transfer or literal translation.

To further explore this trend, we train a series of 1 billion parameter models, each trained on 100 billion tokens of varying proportions of French and English web data. Even at this small scale, we see a trend on at least some standard benchmarks that the pretrained models with at least a 1:2 French to English ratio in pretraining perform better on French versions of standard benchmarks; they also perform significantly better on EIFFEL. This suggests that EIFFEL may serve as an early benchmark for training and that the trend observed with the larger models is already visible at a small scale, meaning that insights from our tests on EIFFEL could inform the training of full-scale LLMs.³

Insofar as idioms are just one example of everyday language likely found in French web data and also an example of a phenomenon important for a variety of downstream tasks—from summarizing conversation transcripts to speaking to users in a style that makes them feel comfortable—our benchmark results show that modulating the amount of French data may be important for downstream success of French LLMs more generally.

An additional factor uncovered by EIFFEL is the effect of translation data. Translation corpora,

¹We take anglocentric models to be those with a high English/French ratio.

²Gallo-centric models have higher proportions of French over English, typically at least 25%.

³We will openly release our entire dataset and our 1B models upon paper acceptance.

ictionaries, etc. imbue models with complex statistical relations between elements of English and French, from the word level up to the sentence, to the paragraph and potentially beyond. EIFFEL shows that a translational paradigm offers rather restricted correlations: models with translation data but not much French in pretraining excel at idioms with word-for-word translations in English, but their performance falls below that of gallo-centric models on other idioms. This supports prior work that has shown that such data induces biases towards frequent or standardized forms over rare and non standard ones, structural simplification as well as reduced lexical and morphological diversity (Vanmassenhove et al., 2021; Laviosa, 1998).

In sum, our main contributions are as follows:

- (i) EIFFEL, a cultural French benchmark for idioms – the only one of its kind.
- (ii) Detailed, reproducible methodology for building such a benchmark.
- (iii) Six 1 billion parameter, open-source models trained from scratch on varying proportions of French and English web data.
- (iv) Experiments on our test models to study the impact of different proportions of French data.

2 State of the Art

The impact of anglocentricity. Anglocentric multilingual models generally can produce reasonable quality non-English text. This does not mean, however, that the concepts and linguistic patterns thereby produced naturally represent concepts and patterns employed by native speakers.

Guo et al. (2025) show that the syntax and vocabulary distribution in non English-languages are affected by high proportions of English training data, resulting in outputs that are often less natural and less diverse than those of native speakers. The greater the typological difference between English and the target language, the more pronounced the gap of lexical naturalness. Karim et al. (2025) show that anglocentricity can impact model performance even in domains that do not seem culturally sensitive, such as math. In particular, they show a decline in performance on mathematical benchmarks when certain words in the math problems are replaced with words more relevant to a non-anglophone culture, such as replacing Western food names with those from Pakistan or Moldova.

Towards multilingual models. Recently, mainstream models including Gemma 3 (Team

et al., 2025), Qwen 3 (Yang et al., 2025), and Mistral 3.1 (<https://mistral.ai/fr/news/mistral-small-3-1>) claim to have increased multilingual data, though we were unable to find statistics on data proportions (and we tried). Some models, such as Llama 3 (Grattafiori et al., 2024), Nemotron-H (Blakeman et al., 2025; Adler et al., 2024), and SmoLLM3 (Bakouch et al., 2025), provide statistics for overall multilingual proportions, but we could not find a breakdown by language.⁴

Projects with a more explicit multilingual focus often provide more information. The EuroLLM models (Martins et al., 2025a,b) cover 24 languages and cite around 45-60% (depending on the training phase) non-English, natural language data with 5-6% in French; the Apertus models (Hernández-Cano et al., 2025) cover 1800 languages and use 40% non-English data with 7.28% French; and the Salamandra models (Gonzalez-Agirre et al., 2025) cover 35 languages and have 55% non-English data with 6.6% in French (and 16% Spanish).

Some projects focus on particular non-English languages; we focus on gallocentric projects here. CroissantLLM is a bilingual 1.3 billion parameter model trained from scratch on a 1:1 French-English ratio (Faysse et al., 2024) while Lucie 7B (Gouvert et al., 2025) and the Gaperon models (Godey et al., 2025) are trained on a roughly 1:2 ratio.

Standard benchmarks. Many benchmarks used to test multilingual performance are translated from datasets originally in English. Some are translated automatically, e.g.: XCODAH and XCSQA (Lin et al., 2021), based on CODAH (Chen et al., 2019) and CSQA (Talmor et al., 2019), as well as ARC (Clark et al., 2018), Hellaswag (Zellers et al., 2019), TruthfulQA (Lin et al., 2022), GSM8K (Cobbe et al., 2021) and MMLU (Hendrycks et al., 2020) translated by (Thellmann et al., 2024). Faysse et al. (2024) translated other benchmarks, including ARC and Hellaswag, into French.

A few benchmarks are multilingual through semi-automatic or manual translation, for example Belebele (Bandarkar et al., 2024), Mintaka (Sen et al., 2022) and Global MMLU (Singh et al., 2025). And only a few benchmarks are originally constructed in the target language, such as FQuAD2.0 (d’Hoffschmidt et al., 2020; Heinrich et al., 2021), a French reading comprehension dataset in the style

⁴Llama 3.1: 8% multilingual data, 8 supported languages; Nemotron-H: 3.7-5%, 9 languages; Nemotron 4 15%, 53 languages; SmoLLM3 12%, 6 languages.

of SQuAD (Rajpurkar et al., 2016).

Benchmarks targeting culture. Global MMLU (Singh et al., 2025) is a multilingual version of MMLU that extends the original benchmark by translating it and tagging questions as culturally-agnostic or culturally-sensitive. BLEnD (Myung et al., 2024) is a multilingual benchmark built by asking native speakers to fill in blanks of translated sentence templates with the names of, say, holidays or common food dishes. CulturalBench (Chiu et al., 2024) includes questions targeting 45 cultures although the questions are written in English.

For cultural benchmarks developed natively in the target language, AraDiCE (Mousi et al., 2025) includes seven Arabic dialogues annotated with associated cultural context. CLiCK (Kim et al., 2024) tests textual, grammatical, and functional knowledge in Korean. IOLBENCH (Goyal and Dan, 2025) poses questions in English about linguistic features of a variety of languages. INCLUDE (Romanou et al., 2024) is a multilingual benchmark built by extracting Q/A data from documents in the target languages that are then verified and corrected by native speakers. It includes culturally agnostic and culturally sensitive subsets. French Bench grammar-vocab-reading (Faysse et al., 2024) evaluates grammar rules, vocabulary, and basic reading comprehension. Of these, only INCLUDE and French Bench cover French, and only INCLUDE has culturally sensitive topics.

Turning to idioms, ID10M (Tedeschi et al., 2022) tests for the ability to identify an idiom or other MWE (multiword expression) in a text. Other tasks focus on being able to provide or identify definitions or paraphrases of idioms and MWEs, such as MAPS (Haviv et al., 2023), IDIOMKB (Li et al., 2024) and MIDAS (Kim et al., 2025). Multilingual Idioms and Similes in LLMs (Khoshtab et al., 2025) tests for the ability to properly continue a text after an idiom is used. Only ID10M includes French, to our knowledge, although the PARSEME shared tasks, e.g., Ramisch et al. (2020), test ability to perform MWE classification and paraphrasing. The closest benchmark to ours is the Arabic benchmark Kinayat (Attia et al., 2025), which assesses the ability to complete idiomatic expressions by masking the last word of the expression.

Ablation studies. There have been few studies that investigate how proportions of data of different languages affect the performance of multilingual models. Han et al. (2025) trained ablation models

on 500 billion Chinese and English tokens in two different settings (1:1 zh-en, 1:9 zh-en). In each setting, they tested replacing both 10 billion and 40 billion regular tokens with the same amount of Chinese-English translation data to evaluate transfer between English and Chinese.

They show that in the 1:9 scenario, replacing 10 billion regular tokens with 10 billion parallel tokens brings Chinese performance to the level of a model trained in the 1:1 setting with no parallel data. This provides experimental evidence that even small amounts of *L1-L2* translation data can improve *L2* performance on standard benchmarks, even if the proportion of *L2* data is small. Han et al. (2025)’s results thus support the multilingual strategy of EuroLLM (Martins et al., 2025a,b). As we show in Section 5, however, this strategy does not capture at least some culturally sensitive aspects of language.

3 Building a benchmark for French idiomatic expressions

The idiomatic meaning of an idiomatic expression cannot be inferred from the literal meanings of its parts: if you’re at a party where no one is talking, and you tell your partner to “break the ice,” you are not literally instructing them to break some block of ice but rather to get people talking. Understanding and properly using idiomatic expressions requires a subtle mastery of the target language and the context of use, making them a perfect subject for a benchmark on culturally-specific language.

EIFFEL draws on the expertise of native speakers. We detail below the steps for building it and illustrate them in Figure 1.

1. Collecting basic idiomatic expressions. Because idiomatic expressions are an important part of everyday language, it was relatively easy to assemble a decent size list by searching the web and discussing among native French and English speaking colleagues. Some expressions were found by starting with English expressions and searching for similar expressions in French.

2. Data categorization. Our hypothesis is that anglocentric multilingual LLMs will fare well on tasks where language transfer helps, but struggle on features that are difficult to capture through translation. Accordingly, we propose three categories of idiomatic expressions for our study:

Word-for-word: The French idiomatic expression has a *word-for-word* translation in English,

e.g., “Ce n’est pas ma tasse de thé” = “It’s not my cup of tea.” We expect these expressions to be the easiest for anglocentric models because they can be inferred from knowledge of the English expression together with basic translation capacities.

Similar: There is an expression in English that is easily recognizable as a translation of the French expression, but is not *word-for-word*, e.g., “appeler un chat un chat” (lit. “call a cat a cat”) vs. “call a spade a spade” or “d’autres chats à fouetter” (lit. “other cats to whip”) vs. “other fish to fry.” We expect anglocentric models to be more likely to confuse the target French expression with a direct French translation of the English expression.

Different: A French expression counts as *different* if we could not find an English counterpart (“de France et de Navarre”) or if the counterpart is sufficiently different that we had to discuss between speakers to find or verify the translations, e.g., “en avoir ras le bol”, which means “have it up to here” but uses the metaphor of a bowl filled to its rim. We hypothesize that these expressions will be the most difficult for models exposed to small percentages of French data.

Of the 602 idiomatic expressions targeted by EIFFEL, 63 are *word-for-word*, 114 are *similar* and 425 are *different*, meaning that EIFFEL emphasizes aspects of language that are particular to French and do not lend themselves to translation.

3. Selection of masking target. As shown in Figure 1, the benchmark is designed as a multiple choice test where the LLM has to fill in a blank (“<...>”) with one of four proposed options to complete the target idiomatic expression. The next step is thus to choose where to put the blank.

This task depends on idiom category. For *word-for-word* expressions, we mask the noun phrase that is most central to the idiom, e.g., “throw the baby out with the bathwater” becomes “throw <...> out with the bathwater”. Note that because French requires adjectives and articles to agree with the head noun for gender (*le* bateau vs. *la* voiture) and number (*les* voitures), these expressions could help the LLM find the correct response. We therefore mask the entire noun phrase.

For *similar* expressions, we aim to mask the most important words that differ between French and English. In general, this involves a noun phrase; “appeler un chat un chat” becomes “appeler <...>”. In rare cases, as when a verb is not widely used in contexts outside the given idiomatic

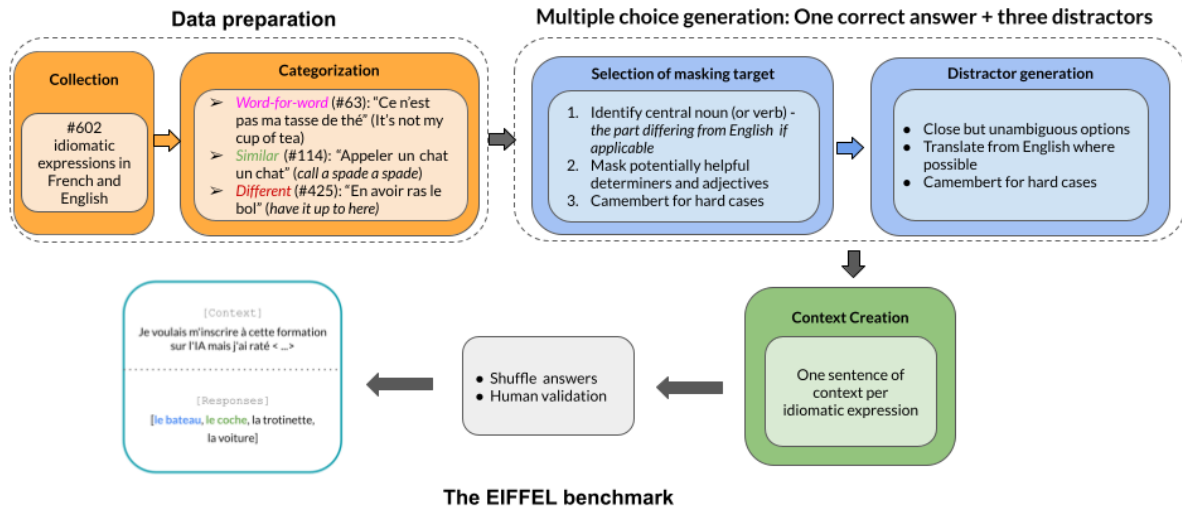


Figure 1: The EIFFEL benchmark building pipeline illustrated here on an example from the *similar* category. The context translates to “I wanted to sign up for this class on AI but I missed <...>”. The possible responses translate as: “the boat,” “the coach,” “the scooter,” “the car”. The correct response is in green; the direct translation of the corresponding English expression is given in blue.

expression or when the verb is the item that differs between the English and French expressions, e.g., “Plonger dans les livres” (lit. “Plunge/dive into the books”) vs. “Hit the books,” we target the verb.

The expressions in the *different* category were more difficult. When we were unable to decide where to put the blank, for any of the categories, we appealed to embeddings by the French model Camembert (Martin et al., 2020). For each alternative under consideration, we looked at the first 15 words whose embeddings were the closest via cosine similarity and chose the alternative with the most pertinent closest neighbors. Impertinent neighbors were those that were close to a non-targeted sense of a polysemous alternative or those whose similarity was not apparent to a native speaker, as can happen for alternatives whose embeddings were clearly not well learned by Camembert. When choosing between two alternatives with pertinent neighbors, we chose the alternative that had the closest neighbors.

4. Distractor generation. Each multiple choice question in the benchmark has one correct answer and three distractors. The latter are crucial for the effectiveness of multiple choice questions and must be both sufficiently credible and unambiguous (Alhazmi et al., 2024).

Given the hypothesis that anglocentric LLMs will have English biases (Guo et al., 2025; Tian et al., 2018), we include the English translation of

masked expressions as distractors when possible, as in “le bateau” (“the boat”) in Figure 1.

When we struggled to choose distractors, we again resorted to Camembert embeddings,⁵ pulling distractors from among the top 15 closest neighbors of the head noun or verb of the masked expression, controlling for grammatical agreement, gender, and semantic compatibility. We avoided neighbors that were so similar that they could lead to answers synonymous with the target expression.

All distractors are human validated for grammar and fluency. To ensure randomness of response order, we shuffled the answers and distractors so that the correct answer is equally likely to show up in any of the four positions.

5. Adding context to idiomatic expressions. Despite our efforts to produce quality, unambiguous distractors, a common issue is that the target sentence could naturally be filled with one or more distractors to make an acceptable French expression. “It’s not my cup of coffee” is a perfect sentence in English (as is its translation in French), but it is not idiomatic. To restrict the task to one of testing for idiomatic expressions, we created contexts for each example that motivated the idiomatic completion. For example: “I don’t like dark chocolate. It’s not

⁵We also tried creating distractors with Mixtral-8x22B-Instruct (<https://huggingface.co/mistralai/Mixtral-8x22B-Instruct-v0.1>) but this approach required significant manual intervention and was generally less satisfying than our Camembert method, so we rejected it.

my cup of <...>.” We constructed and validated all contexts manually. Appendix A provides examples of each of the three categories of EIFFEL.

4 Evaluation of out of the box models

To see whether our benchmark captures performance differences missed by standard benchmarks, we evaluated a series of foundation models on a set of standard benchmarks translated into French and then tested the same models on EIFFEL and the French subset of INCLUDE.

4.1 Models and benchmarks

We restrict our study to base or pretrained models, as we are interested the basic knowledge and linguistic capacities of LLMs. We compare pretrained models in two size ranges, 1-2B and 7-9B, and from three categories: anglocentric, gallo-centric, and intermediate. For anglocentric models, we choose Llama 3.1 8B, Llama 3.2 1B, Gemma 2 9B and Gemma 3 1B. For gallo-centric models, we consider Lucie 7B, Gaperon 1B and 8B, and CroissantLLM (1.3B). As intermediate models, we take EuroLLM 9B and 1.7B, which are trained on less English than anglocentric models, but much less French than more gallo-centric ones.

For benchmarks, we choose a set of standard benchmarks targeting natural language tasks that have both English and French versions: Hellaswag for commonsense reasoning, ARC Challenge for general knowledge and reasoning, MMLU (Global MMLU translations) for general knowledge, and Flores (Costa-Jussà et al., 2022) for translation. For French-centered benchmarks, we consider EIFFEL as well as INCLUDE, whose culturally sensitive subset comes from original French documents (in contrast to Global MMLU for example).

4.2 Evaluation setup

For all of our evaluations, we use the lighteval library (Habib et al., 2023) with the vLLM backend and 0-shot settings. For ARC Challenge, Hellaswag and MMLU, we use normalized accuracy; for FLORES, we use the BLEU metric. When we had an option between cloze or multiple choice formulation, as in MMLU, we chose cloze, which is simpler for pretrained models.

We integrated both INCLUDE and EIFFEL as custom tasks in lighteval with cloze formulation and evaluated on accuracy. For EIFFEL, we consider the log-likelihood of the sequences resulting from completing the context with each response.

4.3 Standard benchmark results

Table 1 shows that all models, even models with substantial French pretraining, tend to do better on the English version of a given benchmark than on its French translation. This tendency is especially marked for anglocentric models Llama3 8B and Gemma 9B with improvements of 7-14 points. We also see a clear improvement with Gaperon 8B, which while trained on the same amount of French data as Lucie 7B, is trained on much more English.

For Lucie and Croissant, the relative improvement on English benchmarks is less pronounced. Given that these models have a 1:1 English French training ratio, switching the language of the benchmark might well have less of an effect.

Perhaps more surprisingly, while anglocentric models tend to be stronger on English versions of the benchmarks than more gallo-centric models, we do not observe the reverse trend in Table 1. Llama and Gemma models tend to have comparable if not slightly better results than Gaperon, Lucie and Croissant on the French benchmarks. Additionally, the EuroLLM models, with an 8:1 to 10:1 English to French ratio, perform more strongly on French benchmarks than the gallo-centric models.

These results suggest several hypotheses. First, given the fact that the French versions of ARC-C, Hellaswag and MMLU are translated from English, one might expect models trained on parallel data to do well on them even if French is not particularly emphasized during training, echoing the results of Han et al. (2025). A second point pertains to the anglocentric orientation of benchmark content: translation should not change the meaning of the original data, so a French version of MMLU will retain the anglocentric biases present in the original dataset. This will give anglocentric models with already high scores on the English versions an advantage even on the French versions.

With regards to Flores, unsurprisingly, the EuroLLM models do very well. The most gallo-centric models Lucie and Croissant are close seconds, however, indicating that a higher French/English ratio does help in translation. More surprisingly, a large majority of the models, including the gallo-centric ones, do better at translating from English to French than vice-versa.

4.4 French-focused benchmarks

The results in Table 1 and our explanatory hypotheses above seem to point to the conclusion that hav-

Pretrained Models	French datasets			English datasets			Translation	
	ARC-C	MMLU	Hellswg	ARC-C	MMLU	Hellswg	En-Fr	Fr-En
Gaperon 8b	.44	.37	<u>.64</u>	.51	.42	.72	.49	<u>.45</u>
Lucie 7b	<u>.40</u>	<u>.35</u>	.65	<u>.39</u>	<u>.41</u>	<u>.67</u>	.51	.47
EuroLLM 9b	.46	.38	.67	.46	<u>.41</u>	.78	.51	.49
Llama-3 8b	.47	.39	.65	.55	.48	.79	<u>.45</u>	.46
Gemma-2 9B	.54	.43	.70	.66	.53	.80	.50	.48
Croissant 1.3b	<u>.28</u>	<u>.28</u>	<u>.50</u>	<u>.27</u>	<u>.31</u>	.53	.45	.42
Gaperon 1.7b	<u>.28</u>	.29	.46	.34	.33	<u>.52</u>	.41	.40
EuroLLM 1.7b	.35	.31	.51	.36	.36	.58	.44	.44
Llama-3 1b	.29	.29	.45	.37	.36	.64	.30	<u>.36</u>
Gemma-3 1b	.30	.30	<u>.50</u>	.38	.36	.62	<u>.26</u>	.39

Table 1: Evaluation of selected models on a set of standard benchmarks with French translations. Models are divided into two categories by size, 1-2 billion parameters and 7-9 billion parameters. High scores are in bold; low scores are underlined. Benchmarks: ARC Challenge (ARC-C), Global MMLU, Hellaswag, FLORES 200 (for translation).

ing high proportions of French data is simply not important for good performance on these datasets. However, another possibility, made more plausible by our hypotheses, is that performing well on the standard benchmarks may not translate to good performance in French on various downstream tasks, requiring, say, conversational fluency.

This possibility motivated us to evaluate our models on the INCLUDE and the EIFFEL benchmarks. As shown in Table 2, less anglocentric models (Gaperon, Lucie, and EuroLLM) tend to outperform more anglocentric models on INCLUDE data judged to be culturally sensitive but not on the culturally agnostic examples. On the sensitive data, Gaperon 8B leads Llama 3 8B by 6 points, while the 1B version comes out 10 points ahead over its Llama counterpart. We note, however, that models with a higher French English data ratio do not always do better on culturally sensitive data; CroissantLLM with 1:1 ratio fares worse than the other gallocentric models with less French data.

On the EIFFEL benchmark, however, overall scores indicate that a higher proportion of French data tends to lead to better performance. When we break down the scores by category, several interesting patterns emerge. We would expect models with a special focus on translation training, such as EuroLLM (Martins et al., 2025b,a) and CroissantLLM (Faysse et al., 2024), to perform well in the *word-for-word* category, and they do. We also expect, however, that models with lower proportions of French data should lose this advantage in the *similar* and *different* categories, where translation is less relevant. Indeed, for these categories,

the gallocentric models beat the EuroLLM models, which in turn beat the anglocentric models.

Test Models	INCLUDE			EIFFEL			
	Ave	Agn	Sens	Ave	W-W	Sim	Diff
Gaperon 8b	.42	.27	.54	.94	<u>.94</u>	.93	.94
Lucie 7b	<u>.37</u>	<u>.23</u>	.51	.94	<u>.94</u>	.94	.94
Eurollm 9b	.39	.25	.51	.93	.98	.88	.92
Llama 3.1 8b	.41	.31	<u>.48</u>	<u>.88</u>	.97	.82	<u>.85</u>
Gemma 9b	.44	.33	.49	.89	.97	<u>.80</u>	.89
Croissant 1.3b	.29	.23	.39	.94	.95	.92	.93
Gaperon 1b	.35	.25	.45	.92	.95	.90	.90
Eurollm 1.7b	.33	.23	.42	.85	.91	.80	.81
Llama 3 1b	<u>.28</u>	<u>.19</u>	<u>.35</u>	<u>.74</u>	<u>.78</u>	.71	<u>.74</u>
Gemma 3 1b	.31	.23	.38	.78	.89	<u>.69</u>	.75

Table 2: Evaluation of selected models on culturally sensitive benchmarks in French. Avg: average, Agn: culturally agnostic, Sens: culturally sensitive, W-W: word for word, Sim: similar, Diff: different.

5 Testing bilingual proportions

Given the difference in performance on standard, translated benchmarks and benchmarks designed for French, we decided to delve deeper into the question of language proportions by training a series of 1 billion parameter models. In addition to monolingual English and French models, we trained four others on varying proportions of French and English web data: 1:100 (fr-en), 1:20, 1:2, and 1:1 (for training details see the Appendix B). The 1:100 and 1:20 mixes represent what we suppose is roughly a minimum and maximum for mainstream LLMs that include, but do not focus on, French. The 1:1 ratio allows for a bilingual

model that is not anglo- (or gallo-) centric, while 1:2 is the proportion of French data used for Lucie and Gaperon. We chose a bilingual approach to limit the number of factors to test and took English as the pivot language, as it plays that role in all mainstream models (as far as we know).

As seen in Table 3, our 1B models fail to provide results significantly beyond a random baseline for French/English ARC challenge except when the models are completely monolingual (French or English). We added ARC-Easy in English⁶ to give the 1B models an easier benchmark, and which with MMLU presents a somewhat less bleak picture. The Hellaswag data set also presents more conclusive results. For these datasets Table 3 shows a performance drop for models below a 1:2 French/English ratio on the French data sets and a corresponding improvement on English datasets once we hit a 1:1 ratio (unfortunately we did not have the resources to test a 2:1 French-English ratio, but we assume the results would be symmetric with the French side).

Fr:En Models	French datasets			English datasets			
	AC	MMLU	HS	AC	AE	MMLU	HS
100% Fr	.26	.26	.38	<u>.20</u>	<u>.32</u>	<u>.25</u>	<u>.28</u>
1:1	.26	.26	.38	.25	.42	.27	.38
1:2	.25	.26	.37	.24	.42	.27	.40
1:20	.25	<u>.25</u>	.32	.24	.44	.27	.42
1:100	<u>.23</u>	<u>.25</u>	.29	.25	.45	.28	.43
0% Fr	.25	<u>.25</u>	<u>.26</u>	.26	.43	.28	.42

Table 3: Evaluation of our 1B models with different ratios of French/English on standard datasets (AC:ARC challenge, AE: ARC Easy, HS: Hellaswag) and their French translations.

Our 1B models’ results on INCLUDE and EIFFEL in Table 4 show that high French to English ratios clearly help with language and culturally sensitive data. Once again below the 1:2 French-English ratio, we see a significant drop in performance on EIFFEL across all categories. INCLUDE reveals a less clear boundary; still, less anglocentric models perform better on the culturally sensitive data than more anglocentric ones.

Performance on EIFFEL, unlike that on INCLUDE, is quite stable and breaks away from random at a low scale. It also includes very high scores and performance improves smoothly throughout training as shown in Figure 5 in the Appendix, indicating that EIFFEL can serve as an early signal benchmark (Penedo et al., 2024). We leave a more

⁶Unfortunately we could not find a French version.

challenging version for future work.

Fr:En Models	INCLUDE			EIFFEL			
	Ave	Agn	Sens	Ave	W-W	Sim	Diff
100% Fr	.27	.19	.35	.89	.89	.89	.89
1:1	.24	.19	.31	.87	.89	.85	.87
1:2	.25	.21	.28	.86	.87	.89	.84
1:20	.28	.23	.33	.70	.79	.68	.63
1:100	<u>.22</u>	<u>.17</u>	<u>.27</u>	.50	.56	.47	.46
0% Fr	.25	.19	.29	<u>.34</u>	<u>.30</u>	<u>.32</u>	<u>.39</u>

Table 4: Results of our 1B models with varying French English training ratios on INCLUDE and EIFFEL.

6 Error analysis on the *similar* category

We did an error analysis of both standard and our 1B models’ performance on the *similar* expressions in EIFFEL. We investigated the total number of errors and looked at how many of these errors resulted from choosing the distractor translated from English as seen in Table 6 in the Appendix. The small anglocentric models Llama, Gemma 1B and our 1B models with ratios of 1:20 Fr-En or less had the most errors (overall and coming from translation) but the lowest *proportion* of literal translation errors. The gallocentric models had lower numbers of overall and translation errors but the number of translation errors varied according to the French/English ratio. Lucie and Croissant with a 1:1 ratio had the lowest number of errors; in Lucie’s case, almost 90% of those errors came from choosing the distractor from English. We also note that our 1B models with a French/English ratio of 1:2 or higher were competitive with EuroLLM 9B. This suggests that a higher French/English training ratio not only improves performance on EIFFEL but allows even the small models to have a fall-back literal translation strategy for difficult idioms.

7 Conclusions

Our experiments with EIFFEL indicate that current multilingual LLMs are often evaluated with tools that insufficiently capture how training data composition shapes model behavior, because of non open data models or lack of testing on multilingual mixes. The dominance of anglocentric resources makes it difficult to disentangle genuine multilingual capabilities from artifacts induced by disproportionate exposure to English. EIFFEL helps correct this imbalance.

658 Limitations

659 Our study focuses on pretrained models, but it
660 would also be relevant to study our question at
661 other stages of model training.

662 A question for our bilingual models that we have
663 not addressed is the question of regional variation.
664 For instance, in Belgian French, the number ninety
665 is expressed as *nonante*, whereas in France, it is
666 *quatre-vingt-dix*. How should we treat these vari-
667 ants is something we leave for future research.

668 Our approach to improving multilingual perfor-
669 mance focuses on a smaller subset of languages,
670 and among high-resource ones rather than attempt-
671 ing to cover all languages simultaneously. This
672 can help to optimize performance for specific lan-
673 guages without diluting training resources. But it is
674 not clear how our approach affects or even transfers
675 to mid or low resource languages. Especially since
676 French and English are close typologically. Our
677 approach also has tested various proportions of the
678 two languages, but absolute counts might matter
679 too. And we have concentrated only on pretraining;
680 different strategies might also be relevant.

681 Another limitation of our work is that the estab-
682 lished benchmarks are not always clean, and even
683 EIFFEL could be improved: many expressions are
684 missing and could be added. We also do not have
685 an English version of EIFFEL to see if the results
686 transfer to English. EIFFEL also takes English
687 as a pivot language, another common assumption
688 but one that may limit the generalizability of the
689 approach.

690 Finally, we did not have the resources to do mul-
691 tiple training runs for our 1B models.

692 References

693 Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H Anh,
694 Pallab Bhattacharya, Annika Brundyn, Jared Casper,
695 Bryan Catanzaro, Sharon Clay, Jonathan Cohen, and
696 1 others. 2024. Nemetron-4 340b technical report.
697 *arXiv preprint arXiv:2406.11704*.

698 Elaf Alhazmi, Quan Z. Sheng, Wei Emma Zhang,
699 Munazza Zaib, and Ahoud Alhazmi. 2024. *Dis-*
700 *tractor generation in multiple-choice tasks: A sur-*
701 *vey of methods, datasets, and evaluation*. *Preprint*,
702 *arXiv:2402.01512*.

703 Mena Attia, Aashiq Muhamed, Mai Alkhamissi,
704 Thamar Solorio, and Mona Diab. 2025. Beyond un-
705 derstanding: Evaluating the pragmatic gap in llms’
706 cultural processing of figurative language. *arXiv*
707 *preprint arXiv:2510.23828*.

Elie Bakouch, Loubna Ben Allal, Anton Lozhkov, Noua-
mane Tazi, Lewis Tunstall, Carlos Miguel Patiño,
Edward Beeching, Aymeric Roucher, Aksel Joonas
Reedi, Quentin Gallouédec, Kashif Rasul, Nathan
Habib, Clémentine Fourier, Hynek Kydlicek, Guil-
herme Penedo, Hugo Larcher, Mathieu Morlon, Vaib-
hav Srivastav, Joshua Lochner, and 4 others. 2025.
SmolLM3: smol, multilingual, long-context reasoner.
<https://huggingface.co/blog/smollm3>.

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel
Artetxe, Satya Narayan Shukla, Donald Husa, Naman
Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and
Madian Khabsa. 2024. The Belebele benchmark: a
parallel reading comprehension dataset in 122 lan-
guage variants. In *Proceedings of the 62nd Annual*
Meeting of the Association for Computational Lin-
guistics (Volume 1: Long Papers), pages 749–775.

Aaron Blakeman, Aarti Basant, Abhinav Khattar,
Adithya Renduchintala, Akhiad Bercovich, Alek-
sander Ficek, Alexis Bjorlin, Ali Taghibakhshi,
Amala Sanjay Deshmukh, Ameya Sunil Mahabalesh-
warkar, and 1 others. 2025. Nemetron-h: A family
of accurate and efficient hybrid mamba-transformer
models. *arXiv preprint arXiv:2504.03624*.

Michael Chen, Mike D’Arcy, Alisa Liu, Jared Fer-
nandez, and Doug Downey. 2019. *Codah: An*
adversarially-authored question answering dataset
for common sense. In *Proceedings of the 3rd Work-*
shop on Evaluating Vector Space Representations for
NLP, pages 63–69, Minneapolis, USA. Association
for Computational Linguistics.

Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin,
Chan Young Park, Shuyue Stella Li, Sahithya Ravi,
Meher Bhatia, Maria Antoniak, Yulia Tsvetkov,
Vered Shwartz, and 1 others. 2024. CulturalBench: a
robust, diverse and challenging benchmark on mea-
suring (the lack of) cultural knowledge of LLMs.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,
Ashish Sabharwal, Carissa Schoenick, and Oyvind
Tafjord. 2018. Think you have solved question an-
swering? try ARC, the AI2 reasoning challenge.
arXiv preprint arXiv:1803.05457.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, and 1 others. 2021. Training verifiers
to solve math word problems. *arXiv preprint*
arXiv:2110.14168.

Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha
Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe
Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,
and 1 others. 2022. No language left behind: Scaling
human-centered machine translation. *arXiv preprint*
arXiv:2207.04672.

Martin d’Hoffschmidt, Wacim Belblidia, Tom Brendlé,
Quentin Heinrich, and Maxime Vidal. 2020. *FQuAD:*
French question answering dataset. *Preprint*,
arXiv:2002.06071.

766	Manuel Faysse, Patrick Fernandes, Nuno M Guerreiro,	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	822
767	António Loison, Duarte Miguel Alves, Caio Corro,	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	823
768	Nicolas Boizard, João Alves, Ricardo Rei, Pedro Hen-	2020. Measuring massive multitask language under-	824
769	rique Martins, and 1 others. 2024. CroissantLLM:	standing. <i>arXiv preprint arXiv:2009.03300</i> .	825
770	A truly bilingual French-English language model.		
771	<i>Transactions on Machine Learning Research</i> .		
772	Nathan Godey, Wissam Antoun, Rian Touchent, Rachel	Alejandro Hernández-Cano, Alexander Hägele,	826
773	Bawden, Éric de la Clergerie, Benoît Sagot, and	Allen Hao Huang, Angelika Romanou, Antoni-Joan	827
774	Djamé Seddah. 2025. Gaperon: A peppered english-	Solergibert, Barna Pasztor, Bettina Messmer, Dhia	828
775	french generative language model suite . <i>Preprint</i> ,	Garbaya, Eduard Frank Ďurech, Ido Hakimi, and	829
776	arXiv:2510.25771.	1 others. 2025. Apertus: Democratizing open and	830
777		compliant llms for global language environments.	831
778	Aitor Gonzalez-Agirre, Marc Pàmies, Joan Llop,	<i>arXiv preprint arXiv:2509.14233</i> .	832
779	Irene Baucells, Severino Da Dalt, Daniel Tamayo,	Aabid Karim, Abdul Karim, Bhoomika Lohana, Matt	833
780	José Javier Saiz, Ferran Espuña, Jaume Prats, Javier	Keon, Jaswinder Singh, and Abdul Sattar. 2025.	834
781	Aula-Blasco, and 1 others. 2025. Salamandra techni-	Lost in cultural translation: Do LLMs struggle with	835
	cal report. <i>arXiv preprint arXiv:2502.08489</i> .	math across cultural contexts? <i>arXiv preprint</i>	836
		<i>arXiv:2503.18018</i> .	837
782	Olivier Gouvert, Julie Hunter, Jérôme Louradour,	Paria Khoshtab, Danial Namazifard, Mostafa Masoudi,	838
783	Christophe Cerisara, Evan Dufraisie, Yaya Sy, Laura	Ali Akhgary, Samin Mahdizadeh Sani, and Yadollah	839
784	Rivière, Jean-Pierre Lorré, and 1 others. 2025.	Yaghoobzadeh. 2025. Comparative study of multilin-	840
785	The Lucie-7b LLM and the Lucie training dataset:	gual idioms and similes in large language models . In	841
786	Open resources for multilingual language generation.	<i>Proceedings of the 31st International Conference on</i>	842
787	<i>arXiv preprint arXiv:2503.12294</i> .	<i>Computational Linguistics</i> , pages 8680–8698, Abu	843
788		Dhabi, UAE. Association for Computational Linguis-	844
789	Satyam Goyal and Soham Dan. 2025. Iolbench: Bench-	tics.	845
790	marking llms on linguistic reasoning. <i>arXiv preprint</i>		
	<i>arXiv:2501.04249</i> .	Eunsu Kim, Juyoung Suk, Philhoon Oh, Haneul Yoo,	846
791	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	James Thorne, and Alice Oh. 2024. CLICk: A bench-	847
792	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	mark dataset of cultural and linguistic intelligence in	848
793	Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,	Korean. <i>arXiv preprint arXiv:2403.06412</i> .	849
794	Alex Vaughan, and 1 others. 2024. The llama 3 herd		
795	of models. <i>arXiv preprint arXiv:2407.21783</i> .	Jisu Kim, Youngwoo Shin, Uji Hwang, Jihun Choi,	850
796		Richeng Xuan, and Taeuk Kim. 2025. Memorization	851
797	Yanzhu Guo, Simone Conia, Zelin Zhou, Min Li, Saloni	or reasoning? exploring the idiom understanding	852
798	Potdar, and Henry Xiao. 2025. Do large language	of llms. In <i>Proceedings of the 2025 Conference on</i>	853
799	models have an english accent? Evaluating and im-	<i>Empirical Methods in Natural Language Processing</i> ,	854
800	proving the naturalness of multilingual LLMs. In	pages 21689–21710.	855
801	<i>Proceedings of the 63rd Annual Meeting of the As-</i>		
802	<i>sociation for Computational Linguistics (Volume 1:</i>	Sara Laviosa. 1998. Core patterns of lexical use in a	856
	<i>Long Papers)</i> , pages 3823–3838.	comparable corpus of english narrative prose. <i>Meta</i> ,	857
803	Nathan Habib, Clémentine Fourrier, Hynek Kydlíček,	43(4):557–570.	858
804	Thomas Wolf, and Lewis Tunstall. 2023. Lighteval:	Shuang Li, Jiangjie Chen, Siyu Yuan, Xinyi Wu, Hao	859
805	A lightweight framework for llm evaluation .	Yang, Shimin Tao, and Yanghua Xiao. 2024. Trans-	860
806		late meanings, not just words: Idiomkb’s role in opti-	861
807	Wenhan Han, Yifan Zhang, Zhixun Chen, Binbin Liu,	mizing idiomatic translation with language models.	862
808	Haobin Lin, Bingni Zhang, Taifeng Wang, Mykola	In <i>Proceedings of the AACL Conference on Artificial</i>	863
809	Pechenizkiy, Meng Fang, and Yin Zheng. 2025.	<i>Intelligence</i> , volume 38, pages 18554–18563.	864
810	Mubench: Assessment of multilingual capabilities of		
811	large language models across 61 languages. <i>arXiv</i>	Bill Yuchen Lin, Seyeon Lee, Xiaoyang Qiao, and Xi-	865
	<i>preprint arXiv:2506.19468</i> .	ang Ren. 2021. Common sense beyond English:	866
812		Evaluating and improving multilingual language	867
813	Adi Haviv, Ido Cohen, Jacob Gidron, Roei Schuster,	models for commonsense reasoning . In <i>Proceedings</i>	868
814	Yoav Goldberg, and Mor Geva. 2023. Understand-	<i>of the 59th Annual Meeting of the Association for</i>	869
815	ing transformer memorization recall through idioms.	<i>Computational Linguistics and the 11th International</i>	870
816	In <i>Proceedings of the 17th Conference of the Euro-</i>	<i>pean Chapter of the Association for Computational</i>	871
817	<i>Linguistics</i> , pages 248–264.	<i>Joint Conference on Natural Language Processing</i>	872
818		(Volume 1: Long Papers), pages 1274–1287, Online.	873
819	Quentin Heinrich, Gautier Viaud, and Wacim Belb-	Association for Computational Linguistics.	
820	lidia. 2021. FQuAD2.0: French question answer-	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	874
821	ing and knowing that you know nothing . <i>Preprint</i> ,	TruthfulQA: Measuring how models mimic human	875
	arXiv:2109.13209.	falsehoods . In <i>Proceedings of the 60th Annual Meet-</i>	876
		<i>ing of the Association for Computational Linguistics</i>	877

878	(<i>Volume 1: Long Papers</i>), pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. <i>arXiv preprint arXiv:1606.05250</i> .	932
879			933
880	Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. 2024. Fineweb-edu: the finest collection of educational content .		934
881			935
882		Carlos Ramisch, Agata Savary, Bruno Guillaume, Jakub Waszczuk, Marie Candito, Ashwini Vaidya, Verginica Barbu Mititelu, Archana Bhatia, Uxoa Iñurieta, Voula Giouli, Tunga Güngör, Menghan Jiang, Timm Lichte, Chaya Liebeskind, Johanna Monti, Renata Ramisch, Sara Stymne, Abigail Walsh, and Hongzhi Xu. 2020. Edition 1.2 of the PARSEME shared task on semi-supervised identification of verbal multiword expressions . In <i>Proceedings of the Joint Workshop on Multiword Expressions and Electronic Lexicons</i> , pages 107–118, online. Association for Computational Linguistics.	936
883	Louis Martin, Benjamin Muller, Pedro Ortiz Suarez, Yoann Dupont, Laurent Romary, Éric Villemonte de La Clergerie, Djamel Seddah, and Benoît Sagot. 2020. Camembert: a tasty french language model. In <i>Proceedings of the 58th annual meeting of the association for computational linguistics</i> , pages 7203–7219.		937
884			938
885			939
886			940
887			941
888			942
889			943
890	Pedro Henrique Martins, João Alves, Patrick Fernandes, Nuno M Guerreiro, Ricardo Rei, Amin Farajian, Mateusz Klimaszewski, Duarte M Alves, José Pombal, Nicolas Boizard, and 1 others. 2025a. Eurollm-9b: Technical report. <i>arXiv preprint arXiv:2506.04079</i> .		944
891			945
892			946
893		Angelika Romanou, Negar Foroutan, Anna Sotnikova, Zeming Chen, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Mohamed A Haggag, Alfonso Amayuelas, and 1 others. 2024. Include: Evaluating multilingual language understanding with regional knowledge. <i>arXiv preprint arXiv:2411.19799</i> .	947
894			948
895			949
896	Pedro Henrique Martins, Patrick Fernandes, João Alves, Nuno M Guerreiro, Ricardo Rei, Duarte M Alves, José Pombal, Amin Farajian, Manuel Faysse, Mateusz Klimaszewski, and 1 others. 2025b. Eurollm: Multilingual language models for europe. <i>Procedia Computer Science</i> , 255:53–62.		950
897			951
898			952
899			953
900			954
901	Bettina Messmer, Vinko Sabolčec, and Martin Jaggi. 2025. Enhancing multilingual llm pretraining with model-based data selection . <i>arXiv</i> .	Priyanka Sen, Alham Fikri Aji, and Amir Saffari. 2022. Mintaka: A complex, natural, and multilingual dataset for end-to-end question answering. <i>arXiv preprint arXiv:2210.01613</i> .	955
902			956
903			957
904			958
905	Basel Mousi, Nadir Durrani, Fatema Ahmad, Md Arid Hasan, Maram Hasanain, Tameem Kabbani, Fahim Dalvi, Shammur Absar Chowdhury, and Firoj Alam. 2025. AraDiCE: Benchmarks for dialectal and cultural capabilities in LLMs. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 4186–4218.	Shivalika Singh, Angelika Romanou, Clémentine Fourier, David Ifeoluwa Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, and 1 others. 2025. Global MMLU: Understanding and addressing cultural and linguistic biases in multilingual evaluation. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 18761–18799.	959
906			960
907			961
908			962
909			963
910			964
911	Junho Myung, Nayeon Lee, Yi Zhou, Jiho Jin, Rifki Putri, Dimosthenis Antypas, Hsuvas Borkakoty, Eunsu Kim, Carla Perez-Almendros, Abinew Ali Ayele, and 1 others. 2024. Blend: A benchmark for LLMs on everyday knowledge in diverse cultures and languages. <i>Advances in Neural Information Processing Systems</i> , 37:78104–78146.	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	965
912			966
913			967
914			968
915			969
916			970
917			971
918	Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. 2024. The fineweb datasets: Decanting the web for the finest text data at scale . In <i>The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. <i>arXiv preprint arXiv:2503.19786</i> .	972
919			973
920			974
921			975
922			976
923			977
924			978
925	Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Amir Hossein Kargaran, Colin Raffel, Martin Jaggi, Leandro Von Werra, and Thomas Wolf. 2025. Fineweb2: One pipeline to scale them all – adapting pre-training data processing to every language . <i>Preprint</i> , arXiv:2506.20920.	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> .	979
926			980
927			981
928			982
929			983
930			984
931			985
			986
			987

988 Simone Tedeschi, Federico Martelli, and Roberto Nav-
989 igli. 2022. Id10m: Idiom identification in 10 lan-
990 guages. In *Findings of the Association for Computa-*
991 *tional linguistics: NAACL 2022*, pages 2715–2726.

992 Klaudia Thellmann, Bernhard Stadler, Michael Fromm,
993 Jasper Schulze Buschhoff, Alex Jude, Fabio Barth,
994 Johannes Leveling, Nicolas Flores-Herr, Joachim
995 Köhler, René Jäkel, and 1 others. 2024. Towards
996 multilingual LLM evaluation for european languages.
997 *arXiv preprint arXiv:2410.08928*.

998 Ye Tian, Ioannis Douratsos, and Isabel Groves. 2018.
999 Treat the system like a human student: Automatic
1000 naturalness evaluation of generated text without refer-
1001 ence texts. In *Proceedings of the 11th International*
1002 *Conference on Natural Language Generation*, pages
1003 109–118.

1004 Eva Vanmassenhove, Dimitar Shterionov, and Matthew
1005 Gwilliam. 2021. Machine translationese: Effects of
1006 algorithmic bias on linguistic complexity in machine
1007 translation. *arXiv preprint arXiv:2102.00287*.

1008 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
1009 Binyuan Hui, Bo Zheng, Bowen Yu, Chang
1010 Gao, Chengen Huang, Chenxu Lv, and 1 others.
1011 2025. Qwen3 technical report. *arXiv preprint*
1012 *arXiv:2505.09388*.

1013 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali
1014 Farhadi, and Yejin Choi. 2019. Hellaswag: Can a
1015 machine really finish your sentence? *arXiv preprint*
1016 *arXiv:1905.07830*.

A Benchmark examples

1017

Target expression	Context	Answer A	Answer B	Answer C	Answer D
Battre le fer quand il est encore chaud <i>Strike while the iron is hot</i>	Tu as eu raison de prendre la parole, il fallait battre < ...> quand il était encore chaud <i>You were right to speak up, we had to strike while < ...> was hot.</i>	le métal <i>the metal</i>	le fer <i>the iron</i>	l'acier <i>the steel</i>	le cuivre <i>the copper</i>
Ce n'est pas ma tasse de thé <i>It's not my cup of tea.</i>	Le chocolat noir ce n'est pas trop ma tasse < ...>. <i>Dark chocolate isn't really my cup <...>.</i>	de thé <i>of tea</i>	d'infusion <i>of infusion</i>	de café <i>of coffee</i>	de tisane <i>of herbal tea</i>
Chercher une aiguille dans une botte de foin <i>Looking for a needle in a haystack</i>	Chercher ce restaurant dans Paris sans GPS c'est comme chercher < ...> dans une botte de foin. <i>Looking for this restaurant in Paris without GPS is like looking for < ...> in a haystack.</i>	une seringue <i>a syringe</i>	une épingle <i>a pin</i>	une aiguille <i>a needle</i>	une ficelle <i>a string</i>
Avoir la tête sur les épaules <i>Have a good head on your shoulders</i>	Il s'agirait d'agir comme un adulte et d'avoir < ...> sur les épaules. <i>It would be a matter of acting like an adult and having < ...> on your shoulders.</i>	le cerveau <i>the brain</i>	la tête <i>the head</i>	le cou <i>the neck</i>	la nuque <i>the back of the neck</i>

Figure 2: Examples of the *word-for-word* category of idiomatic expressions. The correct answers are in blue.

Target expression	Context	Answer A	Answer B	Answer C	Answer D
Avoir un chat dans la gorge <i>To have a cat in the throat</i>	Je suis malade depuis samedi, je suis enrhumé et j'ai < ...> dans la gorge. <i>I've been sick since Saturday, I have a cold and I have < ...> in the throat.</i>	une grenouille <i>a frog</i>	un crapaud <i>a toad</i>	un chien <i>a dog</i>	un chat <i>a cat</i>
Appeler un chat un chat <i>To call a cat a cat</i>	Arrête de prendre des pincettes, au bout d'un moment il faut appeler <...> <i>Stop beating around the bush, at some point you have to call < ...></i>	un chien un chien <i>a dog a dog</i>	une bêche une bêche <i>a spade a spade</i>	un chat un chat <i>a cat a cat</i>	une pelle une pelle <i>a shovel a shovel</i>
Boire comme un templier <i>To drink like a templar</i>	Il a une sacrée descente, il boit comme un < ...> <i>He can really hold his liquor, he drinks like <...></i>	chevalier <i>a knight</i>	templier <i>a templar</i>	dauphin <i>a dolphin</i>	poisson <i>a fish</i>
Être au septième ciel <i>To be in the seventh sky</i>	C'est mon parfum de glace préféré, à chaque fois que j'en mange je suis au < ...> <i>It's my favorite ice cream flavor. Every time I eat it, I'm in <...></i>	septième ciel <i>seventh sky</i>	neuvième nuage <i>ninth cloud</i>	cinquième ciel <i>fifth sky</i>	huitième nuage <i>eighth cloud</i>

Figure 3: Examples of the *similar* category of idiomatic expressions

Target expressions	Context	Answer A	Answer B	Answer C	Answer D
<p>Aller se faire cuire un œuf</p> <p><i>Go fly a kite</i></p>	<p>Il m'agaçait tellement avec ses remarques que je lui ai dit d'aller se faire cuire <...></p> <p><i>He annoyed me so much with his comments that I told him to go to boil an <...></i></p>	<p>un poulet.</p> <p><i>a chicken</i></p>	<p>une soupe.</p> <p><i>a soup</i></p>	<p>un œuf.</p> <p><i>an egg</i></p>	<p>un gâteau.</p> <p><i>a cake</i></p>
<p>Appuyer sur le champignon</p> <p><i>To step on the gas</i></p>	<p>Nous étions déjà en retard, alors il a appuyé sur <...>.</p> <p><i>We were already late, so he pressed <...>.</i></p>	<p>l'aubergine</p> <p><i>the eggplant</i></p>	<p>le champignon</p> <p><i>the mushroom</i></p>	<p>la courgette</p> <p><i>the zucchini</i></p>	<p>la tomate</p> <p><i>the tomato</i></p>
<p>Avaler des couleuvres</p> <p><i>make people believe lies</i></p>	<p>On me fait avaler des <...> toute la journée, répétait le baron.</p> <p><i>They make me swallow <...> all day long, the baron repeated.</i></p>	<p>couleuvres</p> <p><i>grass snakes</i></p>	<p>grenouilles</p> <p><i>frogs</i></p>	<p>lézards</p> <p><i>lizards</i></p>	<p>vipères</p> <p><i>vipers</i></p>
<p>Avoir des oursins dans les poches</p> <p><i>To have deep pockets but short arms.</i></p>	<p>Il refuse toujours de payer un café, ce type a vraiment <...> dans les poches !</p> <p><i>He still refuses to pay for coffee, that guy really has <...> in his pockets!</i></p>	<p>des oursins</p> <p><i>earwigs</i></p>	<p>des poissons</p> <p><i>fish</i></p>	<p>des coquillages</p> <p><i>seashells</i></p>	<p>des épines</p> <p><i>thorns</i></p>

Figure 4: Examples of the *different* category of idiomatic expressions

B Training details for bilingual test models

English data are randomly selected⁷ from the split “sample-350BT” of the FineWeb dataset (Penedo et al., 2024), while French data are taken from the French subset of FineWeb-2 (Penedo et al., 2025). We focus on web data due to their diversity and the fact that they provide the foundation of LLM pretraining. Web data capture a range of everyday language and we assume that they are likely to passively include examples of idiomatic expressions.

We chose the FineWeb and FineWeb-2 datasets because they are relatively recent and have been filtered by similar pipelines. While we could have chosen better filtered datasets, English datasets such as FineWeb-edu (Lozhkov et al., 2024) would have introduced a clear quality difference between the English data and the French data from FineWeb-2. A highly filtered dataset for French, such as FineWeb-2-HQ (Messmer et al., 2025), would have resulted in a dataset too small to train our ablation models on a single epoch.

To tokenize our datasets, we use a custom in-house tokenizer with a vocabulary size of 128,000 that is trained on multilingual data: 20% French, 20% English, 20% Arabic, 20% programming languages and 20% divided between smaller proportions of other European languages. The tokenizer will be made openly available.

Each model is trained on 100 billion tokens and has Llama 3.2 1B architecture except that we adopt a sequence length of 2048 tokens. We use the NeMo library⁸ and the default configuration for the Llama 3.2 1B architecture except that we adopt a sequence length of 2048 tokens. This yields:

number layers	16
hidden size	2048
ffn hidden size	8192
attention heads	32
query groups	8
activation	SwiGLU
normalization	RMS norm

Table 5: Architecture details for our test models.

We adopt a cosine scheduler with a maximum learning rate of $3e^{-4}$ and a minimum learning rate of $3e^{-5}$.

Each model was trained on 64 H100 GPUs (16 nodes) on the Jean Zay supercomputer run

⁷One exception is that we retroactively apply robots.txt protocols, so some samples are excluded off the bat.

⁸<https://github.com/NVIDIA-NeMo/NeMo>.

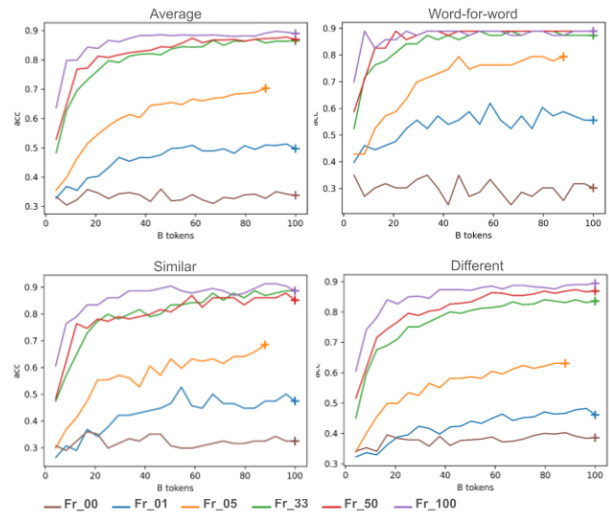


Figure 5: Performance of bilingual test models on the different subsets of EIFFEL. Top left: average scores; top right: performance on *word-for-word* expressions; bottom left: performance on *similar* expression; bottom right: performance on *different* expressions.

by GENCI-IDRIS. Training for each model took around 555 GPU hours.

C Evolution of model training

The four graphs in Figure 5 show how performance of our 1B models evolved on the EIFFEL benchmark throughout training. We see a clear separation between the performance of models with at least a 1:2 ratio of French to English and those with less French. Performance improves fairly smoothly to rise above random performance, indicating that EIFFEL is a good candidate for an early signal benchmark (Penedo et al., 2024).

D Error analysis data

Models below the midline in the Table 6 are our 1B models given in terms of their French:English proportions. The second column in Table 6 provides the total number of errors. The third column indicates the number of errors that result from choosing the distractor coming from the translation of English.

Models	Nb Errors	English Bias
Llama 3 1b	45	23
Gemma 3 1b	35	17
Gemma 9b	23	15
Eurollm 1.7b	23	13
Llama 3.1 8b	20	10
Eurollm 9b	14	10
Gaperon 1b	11	6
Croissant 1.3b	9	6
Gaperon 8b	8	4
Lucie 7b	8	7
0% Fr	77	28
1:100	65	25
1:20	36	14
1:2	13	8
1:1	17	9
100% Fr	13	5

Table 6: Error analysis on the *similar* category of EIF-FEL for all models.